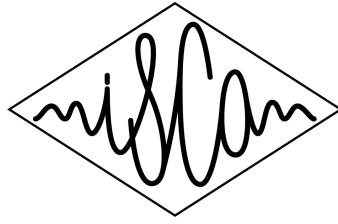# Proceedings of the 8th International Conference on Natural Language and Speech Processing (ICNLSP-2025)

25–27 August, 2025

Order copies of this and other ACL proceedings from:

https://www.icnlsp.org/

# Introduction

It is with great pleasure that we welcome you to the 8th International Conference on Natural Language and Speech Processing (ICNLSP 2025), held at Southern Denmark University from August 25-27, 2025. This volume serves as a comprehensive record of the innovative research and groundbreaking insights on different topics discussed during the conference.

This year's conference attracted 130 submissions from around the globe. The acceptance rate was around 34 %. The 45 accepted papers represent the culmination of rigorous inquiry and intellectual dedication, covering a diverse range of topics within NLP field. Indeed, they showcase the current state of knowledge and shed light on new directions for future exploration. We thank the authors for their valuable contributions.

In order to recognize outstanding scientific contributions, we decided this year to present two awards for the best papers (full and short ones). We congratulate the winners and extend our sincere gratitude to the scientific committee who selected the winners based on recommendations from the program committee, and on the originality, significance, and quality of the research, as well as the clarity of presentation.

We are profoundly honored by the participation of our distinguished keynote speakers, Prof. Dr. Barbara Plank, Prof. Dr. Anders Søgaard, and Prof. Peter Schneider-Kamp. whose insights and vision profoundly enriched the conference.

We thank the conference management members for their efforts, and the program committee and reviewers for their diligent work in curating the high-quality content contained within these pages.

Finally, we are deeply grateful to Southern Denmark University, Danish Data Science Academy, and International Speech Communication Association (ISCA) for their support.

<div align="right">Mourad Abbas, Tariq Yousef, and Lukas Galke</div>

**Organizers:**

**Chairs:**
Dr. Mourad Abbas
Dr. Tariq Yousef

**Program Committee Chair:**
Prof. Lukas Galke

**Publicity Chair:**
Dr. Abed Alhakim Freihat


**Program Committee:**

Ahmed Abdelali, SDAIA, KSA.
Hend Al-Khalifa, KSU, KSA.
Mehmet Fatih Amasyalı, Yildiz Technical University, Turkey.
Yuan An, Drexel University, USA.
Nicklas Sindlev Andersen, SDU, Denmark.
Nicolas Ballier, University Paris Cité, France.
Fayssal Bouarourou, University of Strasbourg, France.
Pierrette Bouillon, University of Geneva, Switzerland.
William Brach, Slovak University of Technology, Slovakia.
Daniel Braun, University of Twente, Netherlands.
Giuseppe Celano, Leipzig University, Germany.
Gérard Chollet, CNRS, France.
Hadda Cherroun, Amar Telidji University, Algeria.
Christian Møller Dahl, SDU, Denmark.
Najim Dehak, Johns Hopkins University, USA.
Alexandra Diehl, University of Zurich, Switzerland.
Andor Diera, Ulm University, Germany.
Dota Dong, Max Planck Institute for Psycholinguistics, Netherlands.
Ashraf Elnagar, University of Sharjah, UAE.
Pascale Feldkamp, Aarhus University, Denmark.
Abed Alhakim Freihat, University of Trento, Italy.
Mounim El Yacoubi, Telecom SudParis, France.
Lukas Galke, SDU, Denmark.
Christian Heumann, Ludwig Maximilian University of Munich, Germany.
Kevin Hirschi, University of Texas at San Antonio, USA.
Torben Johansen, SDU, Denmark.
Kristiina Jokinen, University of Helsinki, Finland.
Eric Laporte, Gustave Eiffel University, France.
Márton Kardos, Aarhus University, Denmark.
Pierre Lison, University of Oslo, Norway.
Mohammed Mediani, UAE University, UAE.

Hermann Ney, RWTH Aachen University, Germany.
Jacob Nielsen, SDU, Denmark.
Martin Rehm, SDU, Denmark.
Uwe Reichel, University of Munich, Germany.
Michael Richter, Leipzig University, Germany.
Shadi Saleh, Prime Technologies, Czech Republic.
Yücel Saygin, Sabanci University, Turkey.
Thomas Schmidt, University of Regensburg, Germany.
Nasredine Semmar, CEA, France.
Rachele Sprugnoli, University of Parma, Italy.
Peter Sullivan, University of British Columbia, Canada.
Irina Temnikova, Big Data for Smart Society Institute, Bulgaria.
María Inés Torres, University of the Basque Country, Spain.
Jan Trmal, AppTek, USA.
Nikos Tsourakis, University of Geneva, Switzerland.
Iraklis Varlamis, Harokopio University of Athens, Greece.
Christian Vedel, SDU, Denmark.
Christian Wartena, University of Applied Sciences and Arts Hannover, Germany.
Ke Yang, University of Texas at San Antonio, USA.
Fayçal Ykhlef, CDTA, Algeria.
Tariq Yousef, SDU, Denmark.
Mai Zaki, American University of Sharjah, UAE.
Zaifu Zhan, University of Minnesota, USA.
Mike Zhang, Aalborg University, Denmark.


**Scientific Committee:**

Prof. Christian Møller Dahl, SDU, Denmark.
Dr. Nicklas Sindlev Andersen, SDU, Denmark.


**Conference Management:**

Lisa Eckford-Soper, SDU, Denmark.
Rula Mreisheh, SDU, Denmark.
Esben Andreas Wrona Bay Sørensen, SDU, Denmark.


**Invited Speakers:**

Prof. Dr. Barbara Plank, LMU Munich, Germany.
Prof. Dr. Anders Søgaard, University of Copenhagen, Denmark.
Prof. Peter Schneider-Kamp, SDU, Denmark.

# Invited Talks

## Human-Centered LLMs for Inclusive Language Technology
*Prof. Dr. Barbara Plank, LMU Munich, Germany*

*Prof. Dr. Barbara Plank* is full professor and chair for AI and Computational Linguistics at LMU Munich, Head of the Munich AI and NLP (MaiNLP) lab, and co-director of the Center for Information and Language Processing (CIS). She is also a visiting full professor at the IT University of Copenhagen.

## What to think of NLP these days?
*Prof. Dr. Anders Søgaard, University of Copenhagen.*

Prof. Dr. Anders Søgaard is a full professor of natural language processing (NLP) and machine learning at the University of Copenhagen in Denmark. Jointly affiliated with the Dpt. of Computer Science, the Dpt. of Philosophy, the Pioneer Centre for Artificial Intelligence, and the Center for Social Data Science. Previously at University of Potsdam, Amazon Core Machine Learning, and Google Research. Father of three and a published poet.

## The Cost of Intelligence: Efficiency Is the Only Path to Democratized AI
*Prof. Peter Schneider-Kamp, SDU, Denmark*

Prof. Peter Schneider-Kamp is a Professor of Computer Science at the University of Southern Denmark (SDU), where he holds a chair in Artificial Intelligence (AI) within the Center for Machine Learning. He co-leads the Danish Foundation Models (DFM) project on multilingual language models with a focus on Danish.

# Table of Contents

# Zero-Shot Commonsense Validation and Reasoning with Large Language Models: An Evaluation on SemEval-2020 Task 4 Dataset

**Rawand Alfugaha**
College of Information Technology
Lusail University
Doha, Qatar
ralfoqha@lu.edu.qa

**Mohammad AL-Smadi**
Digital Learning and Online Education Office
Qatar University
Doha, Qatar
malsmadi@qu.edu.qa

## Abstract

This study evaluates the performance of Large Language Models (LLMs) on SemEval-2020 Task 4 dataset, focusing on commonsense validation and explanation. Our methodology involves evaluating multiple LLMs, including LLaMA3-70B, Gemma2-9B, and Mixtral-8x7B, using zero-shot prompting techniques. The models are tested on two tasks: Task A (Commonsense Validation), where models determine whether a statement aligns with commonsense knowledge, and Task B (Commonsense Explanation), where models identify the reasoning behind implausible statements. Performance is assessed based on accuracy, and results are compared to fine-tuned transformer-based models. The results indicate that larger models outperform previous models and perform closely to human evaluation for Task A, with LLaMA3-70B achieving the highest accuracy of **98.40%** in Task A whereas, lagging behind previous models with **93.40%** in Task B. However, while models effectively identify implausible statements, they face challenges in selecting the most relevant explanation, highlighting limitations in causal and inferential reasoning.

## 1 Introduction

Commonsense reasoning is a crucial aspect of Natural Language Processing (NLP) that enables models to understand and validate knowledge beyond explicit textual data. The motivation behind this research comes from the need to develop NLP models that can reason beyond surface-level text representations and apply real-world knowledge to language understanding tasks. Existing benchmarks, such as CommonGen (Lin et al., 2019), SemEval-2020 Task 4: Commonsense Validation and Explanation (Wang et al., 2020), Common-SenseQA 2.0 (Talmor et al., 2022), and COPEN (Peng et al., 2022), have highlighted various aspects of commonsense reasoning, including generative

commonsense reasoning, multi-hop reasoning, and physical commonsense knowledge. However, these tasks still pose challenges in handling nuanced reasoning (El-Sayed and Pacholczyk, 2002), causal inference(Yao et al., 2021), and knowledge integration (Chen et al., 2020).

The SemEval-2020 Task 4: Commonsense Validation and Explanation (Wang et al., 2020) has served as a benchmark for evaluating various NLP models' capabilities in this domain. The task consistes of three sub-tasks, where in this research we will focus on the first two namely: Task A - Commonsense Validation: Determining whether a given statement aligns with commonsense knowledge, and Task B - Commonsense Explanation: Identifying the reasoning behind why a statement is against common sense. Table 1 provides examples on both tasks as they appear in the dataset.

This paper aims to explore how well large language models (LLMs) perform on commonsense reasoning tasks using zero-shot prompting. By evaluating multiple LLMs on SemEval-2020 Task 4, we investigate their ability to reason effectively without explicit fine-tuning. We present an overview of existing research, detail our methodology, and analyze experimental results to assess the strengths and limitations of current approaches.

## 2 Related Work

SemEval-2020 Task 4, which focuses on Commonsense Validation and Explanation, attracted considerable attention, with numerous teams participating in its three subtasks. This literature review highlights the best-performing models in Tasks A and B, showcasing their methodologies and contributions to the field.

CN-HIT-IT.NLP (Zhang et al., 2020) emerged as the leading model in Subtask A, employing a variant of K-BERT (Liu et al., 2019a) as its encoder. This model stands out for its integration of

| Task | Example |
|------|---------|
| **Task A: Commonsense Validation** | **Which statement is against common sense?** |
| | - **Statement 1:** He put a turkey into the fridge. ( Correct) |
| | - **Statement 2:** He put an elephant into the fridge. (Against commonsense) |
| **Task B: Commonsense Explanation** | **Why is this statement against common sense?** |
| | **Statement:** He put an elephant into the fridge. |
| | - **A:** An elephant is much bigger than a fridge. ( Correct) |
| | - **B:** Elephants are usually white while fridges are usually white. |
| | - **C:** An elephant cannot eat a fridge. |

Table 1: Examples of Commonsense Validation and Explanation Tasks

knowledge graphs, specifically ConceptNet (Speer et al., 2017), which allows it to extract relevant triples that enhance the understanding of language representations. This approach underscores the importance of leveraging structured knowledge to improve commonsense reasoning capabilities.

In Subtask B, ECNU-SenseMaker (Zhao et al., 2020) achieved top performance by also utilizing K-BERT (Liu et al., 2019a). This model innovatively combines structured knowledge from ConceptNet (Speer et al., 2017) with unstructured text through a Knowledge-enhanced Graph Attention Network. This integration facilitates a deeper understanding of commonsense knowledge, demonstrating the effectiveness of combining different types of information to enhance model performance.

Another notable model, IIE-NLP-NUT (Xing et al., 2020), utilized RoBERTa as its encoder. This model's unique contribution lies in its second pretraining phase, which involved a textual corpus from the Open Mind Common Sense (OMCS) project (Singh et al., 2002). By exploring various prompt templates for input construction, this model illustrates the potential of tailored input strategies in improving commonsense validation tasks

Team Solomon (Srivastava et al., 2020) was ranked 5th and 4th in Subtasks A and B, respectively. Their approach, which also relied on RoBERTa, highlighted the capacity of large-scale pretrained language models to encapsulate commonsense knowledge effectively without external resources.

Across the two subtasks, the dominant trend was the use of large-scale pretrained language models such as K-BERT (Liu et al., 2019a), *RoBERTa* (Liu et al., 2019b), *BERT* (Devlin et al., 2018), and *GPT-2* (Radford et al., 2019), often fine-tuned

with additional commonsense knowledge sources. Additionally, models incorporating external structured knowledge sources (e.g., ConceptNet) generally outperformed purely language-model-based approaches.

## 3 Methodology

Our study aims at evaluating the performance of multiple Large Language Models (LLMs) for commonsense validation and reasoning using zero-shot prompting. This approach leverages pre-trained LLMs without task-specific fine-tuning, relying solely on their inherent reasoning capabilities. For this purposes, we utilize the SemEval-2020 Task 4 dataset (Wang et al., 2020), which comprises labeled statements designed for commonsense validation and explanation tasks. To ensure a fair comparison between explicitly fine-tuned models and those evaluated solely with zero-shot prompting, we use only the test set for evaluation. The test set contains 1,000 entries for each task (Task A and Task B), providing a standardized benchmark for assessing model performance. The dataset is publicly available and can be accessed at [1].

As depicted in Figure 1, the methodology consists of the following key stages:

- **Pre-processing:** preparing the input test data templatic prompt to ensure compatibility with zero-shot prompting.

- **Model Calling:** Applying zero-shot prompting to multiple LLMs, including LLaMA3, Gemma2, and Mixtral to assess their commonsense validation and reasoning abilities.

---

[1] https://github.com/wangcunxiang/
SemEval2020-Task4-Commonsense-Validation-and-Explanation

Figure 1: The architecture of the commonsense validation and reasoning with zero-shot prompting of LLMs.

LLMs are directly accessible through the Gro-qCloud [2] Models API endpoint using the model IDs

- **Performance Metrics:** Evaluating model outputs based on accuracy to quantify their effectiveness.

- **Comparative Analysis:** Benchmarking zero-shot LLMs performance against fine-tuned transformer models to examine the impact of training on commonsense validation and reasoning tasks.

## 4   Results and Discussion

Table 2 presents the performance of the models on the commonsense validation (Task A) and commonsense explanation (Task B) tasks from SemEval-2020 Task 4. The results for human performance and transformer-based models (CN-HIT-IT.NLP, ECNU-SenseMaker, IIE-NLP-NUT, and Solomon) are as reported in the original SemEval-2020 Task 4 paper (Wang et al., 2020). In contrast, the results for the LLMs (LLaMA3, Gemma2, and Mixtral) are obtained from our experiments with zero-shot prompting. Findings are reported in the following subsections.

### 4.1   Performance Analysis

Among the models evaluated in this study, **L3-70B (LLaMA3-70B)** demonstrated the highest performance in Task A, scoring **98.4%**, with an evidence that large-scale LLMs can effectively validate commonsense knowledge with zero-shot prompting. However, its performance in Task B (**93.4%**) lags behind the transformer-based models reported as top 4 performing models in the Task paper. These models were explicitly fine-tuned for the task and some of them used external resources for the models training. This indicates that while LLMs are highly proficient in identifying implausible statements, they still struggle with selecting the most relevant explanation, demonstrating limitations in causal and inferential reasoning.

Similarly, the **G2-9B (Gemma2-9B)** model achieves strong performance in Task A (**97.9%**) but showing a more significant decline in Task B (**91.0%**). This further highlights the challenge of explanation selection, as these models may recognize implausibility without fully understanding the underlying causal mechanisms.

A size-dependent trend is observed in the LLaMA3 models. The smaller **L3-8B (LLaMA3-8B)** demonstrates significantly weaker performance than its larger version, with **84.4%** in Task A and **83.1%** in Task B. Finaly, the **M8x7B (Mixtral-8x7B)** model exhibited the weakest performance, with **66.0%** in Task A and **50.9%** in Task B. Its near-random performance in explanation selection

3

| Model | Task A (Validation) (%) | Task B (Explanation) (%) |
|---|---|---|
| Human | 99.1 | 97.8 |
| CN-HIT-IT.NLP | 97.0 | **94.80** |
| ECNU-SenseMaker | 96.7 | 95.0 |
| IIE-NLP-NUT | 96.4 | 94.3 |
| Solomon | 96.0 | 94.0 |
| L3-70B (LLaMA3-70B) | **98.40** | 93.40 |
| G2-9B (Gemma2-9B) | 97.90 | 91.00 |
| L3-8B (LLaMA3-8B) | 84.40 | 83.10 |
| M8x7B (Mixtral-8x7B) | 66.00 | 50.90 |

Table 2: Performance of different models on Task A (Commonsense Validation) and Task B (Commonsense Explanation) for English data. The models are: L3-70B (LLaMA3-70B), G2-9B (Gemma2-9B), L3-8B (LLaMA3-8B), and M8x7B (Mixtral-8x7B).

| id | sent0 | sent1 | L3-70B | G2-9B | M8x7B | L3-8B |
|---|---|---|---|---|---|---|
| 459 | The dog pounced on the rabbit | The cat pounced on the rabbit | sent0 | sent0 | **Other** | sent0 |
| 737 | She purchased four supermarket tickets. | She purchased four theater tickets. | sent1 | sent1 | sent1 | sent1 |
| 174 | Witches are not made of wood | Toads are not made of wood | sent0 | sent0 | sent0 | sent0 |

Table 3: Sample of common misclassified instances for TaskA. Model abbreviations: L3-70B = LLaMA3-70B, G2-9B = Gemma2-9B, M8x7B = Mixtral-8x7B, L3-8B = LLaMA3-8B. Keep in mind that Task A is about identifying which statement is against common sense?

suggests that it struggles not only with causal inference but also with general commonsense understanding, likely due to limitations in its training data or architecture. It is important to note that this lower accuracy was not due to weak reasoning abilities but rather due to inconsistencies in the output format, where the model provided both classification and explanation instead of following the expected template for the output.

## 4.2 Implications for Zero-Shot Commonsense Reasoning

The results indicate that while LLMs often recognize implausible statements but fail to select the most relevant explanation, highlighting deficits in causal and inferential reasoning. This suggests that current zero-shot approaches may capture surface-level plausibility but lack deeper reasoning abilities necessary for explanation generation.

Furthermore, the comparison between pre-trained LLMs and task-specific models from SemEval-2020 Task 4 suggests that explicit fine-tuning on commonsense explanation data remains beneficial. While larger models such as **L3-70B** outperform fine-tuned models in validation, they do not surpass them in explanation selection, reinforcing the need for additional adaptation to improve causal reasoning.

## 4.3 Common Misclassification Patterns

An analysis of misclassified instances provides insights into the reasoning patterns of different models. In **Task A**, some models failed to differentiate between subtle variations in sentence structure. For example, the model incorrectly classified the following pair:

*The dog pounced on the rabbit. The cat pounced on the rabbit.*

This type of error suggests that the models may rely on statistical patterns rather than deep semantic understanding.

In **Task B**, errors were primarily related to the selection of the most plausible explanation. A notable example is:

**False Statement:** "There are four years each season."
**Correct Explanation:** "A year can be divided into four seasons."

| id | FalseSent. | OptionA | OptionB | OptionC | L3-70B | G2-9B | M8x7B | L3-8B |
|---|---|---|---|---|---|---|---|---|
| 1388 | Roberts' room is sleeping | A room cannot close his eyes, because he has no eyes. | Robert won't let the room sleep because he needs rest. | Robert can sleep in his room | A | A | A | A |
| 1444 | There are four years each season. | Different seasons have different temperatures. | A year can be divided into four seasons. | A season is shorter than a year. | C | C | **Other** | C |
| 1172 | People can need sleep. | Sleep is not a thing to have it granted. | Sleeping is nature for every living being. | Sleeping is an activity that every living thing does. | A | A | C | A |

Table 4: Sample of common misclassified instances for TaskB. Model abbreviations: L3-70B = LLaMA3-70B, G2-9B = Gemma2-9B, M8x7B = Mixtral-8x7B, L3-8B = LLaMA3-8B. Task B is about selecting the reason for Why is this statement against common sense?

Some models selected incorrect explanations, indicating potential limitations in their ability to link cause-effect relationships effectively. It should be noted that sentence IDs **1388, 1444, and 1172** are not present in the common misclassified instances of Task A.

Despite the overall strong performance, the results also highlight challenges in certain reasoning aspects. The models demonstrated difficulty in selecting the most appropriate explanation for an implausible statement in Task B, even though they performed well in identifying implausible statements in Task A. This suggests that while the models recognize commonsense inconsistencies, they may struggle to justify their choices accurately. One possible explanation for this challenge is that Task B requires models to establish causal or inferential relationships between a false statement and its explanation. While Task A is a binary classification task requiring identification of implausible statements, Task B introduces additional complexity by demanding a deeper understanding of reasoning patterns and cause-effect relationships. Selecting the correct explanation requires not only recognizing a logical inconsistency but also evaluating multiple plausible justifications and determining which one best aligns with human commonsense knowledge. This suggests that current LLMs, despite their powerful language modeling capabilities, may still struggle with selecting the most contextually relevant explanation among multiple plausible options, as this task requires a nuanced understanding of real-world implications and reasoning structures (Mondorf and Plank, 2024).

Additionally, the low measured performance of **Mixtral-8x7B** can be attributed to output inconsistencies. The model frequently produced both an answer and an explanation, which deviated from the required response format. This indicates that we cannot rely on the achieved results for this model to evaluate its performance on both tasks. More post-processing steps are required to ensure consistent output formatting when evaluating model performance.

### 4.4 Conclusion

This study demonstrates that large-scale LLMs, particularly **LLaMA3-70B** and **Gemma2-9B**, exhibit strong commonsense reasoning capabilities even in a zero-shot setting. These models outperform state-of-the-art fine-tuned transformer-based models, indicating that LLMs can generalize well across commonsense validation tasks. However, challenges remain in explanation selection and maintaining consistent output formats. Future research may include exploring Commonsense knowledge-graph

LLMs (Li et al., 2022; Zhao et al., 2024; Toroghi et al., 2024), in addition to fine-tuning strategies, retrieval-augmented approaches, and structured prompting techniques to enhance the inferential reasoning capabilities of LLMs in zero-shot settings.

# References

Xiaojun Chen, Shengbin Jia, and Yang Xiang. 2020. A review: Knowledge reasoning over knowledge graph. *Expert systems with applications*, 141:112948.

Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. 2018. Bert: Pre-training of deep bidirectional transformers for language understanding. In *arXiv preprint arXiv:1810.04805*.

Mazen El-Sayed and Daniel Pacholczyk. 2002. A qualitative reasoning with nuanced information. In *Logics in Artificial Intelligence*, pages 283–295, Berlin, Heidelberg. Springer Berlin Heidelberg.

Xiang Lorraine Li, Adhiguna Kuncoro, Jordan Hoffmann, Cyprien de Masson d'Autume, Phil Blunsom, and Aida Nematzadeh. 2022. A systematic investigation of commonsense knowledge in large language models. In *Proceedings of the 2022 Conference on Empirical Methods in Natural Language Processing*, pages 11838–11855, Abu Dhabi, United Arab Emirates. Association for Computational Linguistics.

Bill Yuchen Lin, Wangchunshu Zhou, Ming Shen, Pei Zhou, Chandra Bhagavatula, Yejin Choi, and Xiang Ren. 2019. Commongen: A constrained text generation challenge for generative commonsense reasoning.

Weijie Liu, Peng Zhou, Zhe Zhao, Zhiruo Wang, Qi Ju, Haotang Deng, and Ping Wang. 2019a. K-bert: Enabling language representation with knowledge graph. In *arXiv preprint arXiv:1909.07606*.

Yinhan Liu, Myle Ott, Naman Goyal, Jingfei Du, Mandar Joshi, Danqi Chen, Omer Levy, Mike Lewis, Luke Zettlemoyer, and Veselin Stoyanov. 2019b. Roberta: A robustly optimized bert pretraining approach. In *arXiv preprint arXiv:1907.11692*.

Philipp Mondorf and Barbara Plank. 2024. Beyond accuracy: Evaluating the reasoning behavior of large language models–a survey. *arXiv preprint arXiv:2404.01869*.

Hao Peng, Xiaozhi Wang, Shengding Hu, Hailong Jin, Lei Hou, Juanzi Li, Zhiyuan Liu, and Qun Liu. 2022. Copen: Probing conceptual knowledge in pre-trained language models. *arXiv preprint arXiv:2211.04079*.

Alec Radford, Jeff Wu, Rewon Child, David Luan, Dario Amodei, and Ilya Sutskever. 2019. Language models are unsupervised multitask learners. In *Technical Report*.

Push Singh, Thomas Lin, Erik T. Mueller, Grace Lim, Travell Perkins, and Wan Li Zhu. 2002. Open mind common sense: Knowledge acquisition from the general public. In *On the Move to Meaningful Internet Systems 2002: CoopIS, DOA, and ODBASE*, pages 1223–1237, Berlin, Heidelberg. Springer Berlin Heidelberg.

Robyn Speer, Joshua Chin, and Catherine Havasi. 2017. Conceptnet 5.5: an open multilingual graph of general knowledge. In *Proceedings of the Thirty-First AAAI Conference on Artificial Intelligence*, AAAI'17, page 4444–4451. AAAI Press.

Vertika Srivastava, Sudeep Kumar Sahoo, Yeon Hyang Kim, Rohit R.r, Mayank Raj, and Ajay Jaiswal. 2020. Team Solomon at SemEval-2020 task 4: Be reasonable: Exploiting large-scale language models for commonsense reasoning. In *Proceedings of the Fourteenth Workshop on Semantic Evaluation*, pages 585–593, Barcelona (online). International Committee for Computational Linguistics.

Alon Talmor, Ori Yoran, Ronan Le Bras, Chandra Bhagavatula, Yoav Goldberg, Yejin Choi, and Jonathan Berant. 2022. Commonsenseqa 2.0: Exposing the limits of ai through gamification.

Armin Toroghi, Willis Guo, Mohammad Mahdi Abdollah Pour, and Scott Sanner. 2024. Right for right reasons: Large language models for verifiable commonsense knowledge graph question answering. In *Proceedings of the 2024 Conference on Empirical Methods in Natural Language Processing*, pages 6601–6633, Miami, Florida, USA. Association for Computational Linguistics.

Cunxiang Wang, Shuailong Liang, Yili Jin, Yilong Wang, Xiaodan Zhu, and Yue Zhang. 2020. Semeval-2020 task 4: Commonsense validation and explanation. In *Proceedings of The 14th International Workshop on Semantic Evaluation*. Association for Computational Linguistics.

Luxi Xing, Yuqiang Xie, Yue Hu, and Wei Peng. 2020. IIE-NLP-NUT at SemEval-2020 task 4: Guiding PLM with prompt template reconstruction strategy for ComVE. In *Proceedings of the Fourteenth Workshop on Semantic Evaluation*, pages 346–353, Barcelona (online). International Committee for Computational Linguistics.

Liuyi Yao, Zhixuan Chu, Sheng Li, Yaliang Li, Jing Gao, and Aidong Zhang. 2021. A survey on causal inference. *ACM Transactions on Knowledge Discovery from Data (TKDD)*, 15(5):1–46.

Yice Zhang, Jiaxuan Lin, Yang Fan, Peng Jin, Yuanchao Liu, and Bingquan Liu. 2020. Cn-hit-it.nlp at semeval-2020 task 4: Enhanced language representation with multiple knowledge triples. In *Proceedings of The 14th International Workshop on Semantic Evaluation*. Association for Computational Linguistics.

Qian Zhao, Siyu Tao, Jie Zhou, Linlin Wang, and Xin Lin. 2020. Ecnu-sensemaker at semeval-2020 task 4: Leveraging heterogeneous knowledge resources for commonsense validation and explanation. In *Proceedings of The 14th International Workshop on Semantic Evaluation*. Association for Computational Linguistics.

Zirui Zhao, Wee Sun Lee, and David Hsu. 2024. Large language models as commonsense knowledge for large-scale task planning. *Advances in Neural Information Processing Systems*, 36.

# Bridging the Gap: Design and Evaluation of an Automated System for French Cued Speech

**Brigitte Bigi**

Laboratoire Parole et Langage, CNRS, Aix Marseille Univ, 13100 Aix-en-Provence - France
`brigitte.bigi@cnrs.fr`

## Abstract

Access to spoken language remains a challenge for deaf and hard-of-hearing individuals due to the limitations of lipreading. Cued Speech (CS) addresses this by combining lip movements with hand cues—specific shapes and placements near the face—making each syllable visually distinct. This system complements cochlear implants and supports oral language, phonological awareness, and literacy. This paper introduces the first open-source system for automatically generating CS in video format. It takes as input a video recording, the corresponding audio signal, and an orthographic transcript. These elements are processed through a modular pipeline, which includes phonetic mapping, temporal alignment, spatial placement, and real-time rendering of a virtual coding hand. The system is multilingual by design, with current resources focused on French. An evaluation under varied conditions showed decoding rates up to 92% for manually coded stimuli, and averages exceeding 80% for automatically generated ones. Visual clarity of hand shapes proved more critical than timing or angle. Stylized designs and frontal views enhanced decoding performance, while attempts at naturalistic rendering or visual effects often degraded it. These findings indicate that visual abstraction improves readability. This work provides a reproducible and scientifically grounded framework for visual phonetic encoding, and delivers a practical tool for education, accessibility, and research.

## 1 Introduction

### 1.1 Visual Access to Spoken Language through Cued Speech

Speech production involves both acoustic and visual cues. While lip movements convey useful information, many phonemes appear identical on the lips and form so-called "visemes"—groups of phonemes that are visually indistinguishable (Fisher, 1968; Massaro and Palmer Jr, 1998). As a result, lipreading remains highly ambiguous: correct word identification rarely exceeds 30% (Rönnberg, 1995; Rönnberg et al., 1998).

To address this limitation, R. Orin Cornett introduced Cued Speech (CS) (Cornett, 1967), a visual communication system designed to make each phoneme visually distinct. CS combines lip movements with hand cues—specific handshapes and positions placed around the face—that encode consonants and vowels. It provides full visual access to spoken language and supports phonological awareness, literacy development, and spoken language acquisition in deaf or hard-of-hearing individuals (Clarke and Ling, 1976; Neef and Iwata, 1985). CS has since been adapted to over 65 languages[1].

Cued Speech is widely used by speech-language pathologists to support early language acquisition in deaf children. Among others, in France, it is promoted by the Association pour la Langue française Parlée Complétée (ALPC)[2], and in the US by the National Cued Speech Association[3]. Numerous studies have shown that CS enhances access to phonological structure, supports literacy development, and fosters inclusive education (Leybaert and Charlier, 1996; Colin et al., 2017; LaSasso et al., 2010). Together, these findings highlight its importance in supporting language acquisition pathways for deaf learners.

Building on its demonstrated benefits for access to spoken language, Cued Speech and Sign Languages serve distinct linguistic and cultural functions. They are not mutually exclusive: while some deaf children follow a sign language pathway, access to reading and writing typically requires exposure to spoken language. By offering a precise visual representation of sounds, Cued Speech supports this process. It is therefore relevant to all deaf learners aiming to acquire spoken language,

---

[1] `https://www.academieinternationale.org/`
[2] `https://alpc.asso.fr/`
[3] `https://cuedspeech.org/`

whether or not they use a sign language. This distinction is essential to avoid misinterpretations: Cued Speech is not a language and is not intended to replace natural sign languages such as LSF, but to complement them when access to spoken language is required or preferred.

Following the general principles of CS, the French adaptation was developed in the 1970s. It uses eight handshapes to encode consonants and five facial positions to encode vowels. Each syllable is represented by a combination of lip movement and a hand cue, also called a key, formed by a handshape–position pair. A simple syllable like CV or V is coded by a single key, while more complex structures, such as CCV, require multiple successive keys: for example, a 'C-' followed by a 'CV' structure. To illustrate this system, Figure 1 shows the handshapes used for consonants, and Figure 2 shows the vowel positions around the face. Both figures include the neutral position used when no speech is pronounced.



Figure 1: Handshapes representing consonants



Figure 2: Positions representing vowels

Below is a concrete example showing how a sentence is encoded into cues:

```
text:            Tu es gris.
phones:          t y e g ʁ i
CV sequence:     C V V C C V
cues-structure:  C-V.-V.C-.C-V
cues code:       5-t.5-t.7-s.3-m
```

The internal consistency of Cued Speech makes it well-suited for automation. Generating cues from speech or text opens the door to a wide range of applications: cued videos for learning and access, training tools for families and educators, and greater availability of CS in contexts where trained coders are not present. More broadly, automatic cueing can support language acquisition in deaf children, improve communication in mixed hearing environments, and reinforce lipreading skills.

This paper presents the first complete and shareable system for automatic CS generation. It takes as input a video recording, its audio signal, and a transcript, and produces a new version of the video in which a synchronized virtual hand encodes the CS transcription. The architecture was built entirely from scratch, formalizing each stage of the process from segmentation to cue rendering. It is designed for multilingual use and has been implemented and tested for French. The full system is open-source, and all components have been evaluated with end-user testing.

## 1.2 Related Works

The first attempt to automate cue generation, AutoCuer, was developed by R. Orin Cornett himself (Cornett et al., 1977). Between 1995 and 2000, a series of studies at the Massachusetts Institute of Technology (MIT) explored real-time automatic cueing for American English (Bratakos, 1995; Sexton, 1997; Bratakos et al., 1998; Duchnowski et al., 2000). These remain the most advanced documented efforts in the field. Their system relied on speaker-dependent automatic speech recognition to extract phonemes from live recordings, which were then converted into hand cues and displayed as a virtual hand overlaid on the video. Evaluations showed significant gains in decoding accuracy, with some conditions yielding scores twice as high as lipreading alone. However, many components required manual adjustments (Sexton, 1997): cue positions were initialized by hand, transitions were interpolated without formal modeling, and the mapping rules were not described in reusable form. The lack of published code or reproducible design has prevented further development or reuse.

9

To date, no operational or open-source tool exists for automatic CS generation in any language, despite increasing scientific interest and documented benefits.

## 1.3 Foundations and Scope

Developing a complete system was a necessary step, independently of data availability. It provided the opportunity to define a structured architecture, implement a fully functional version, and formalize the modeling of each component. The resulting system is transparent, reproducible, and already usable in real conditions. It operates with minimal computational cost, can be refined through expert feedback, and offers a solid basis for future improvements, including data-driven modules once annotated resources become available.

A French Cued Speech corpus has recently been collected and partially annotated (Bigi et al., 2022), but the annotation process is still ongoing due to the precision required.

This work then marks the beginning of a long-term effort to build a reliable and extensible framework for automatic CS generation. It defines a shared foundation for future developments in augmented video production and evaluation.

## 2 System Description

While many studies describe individual aspects of CS production—such as articulation, speech coordination, timing, or spatial organization—formal descriptions remain rare. Few are presented in a way that supports computational modeling or system implementation. The literature describes many aspects of CS production. However, formal accounts of its speech coordination, timing, and spatial organization remain rare. Few works address these questions, and the descriptions are rarely framed in terms of computational modeling.

In this work, the cueing process was analyzed by combining published linguistic descriptions (Attina, 2005; Aboutabit, 2007) with structured discussions conducted with experienced coders. This led to the identification of four core processing components, which structure the system: determining what to code (i.e., the sequence of keys from the phoneme transcription), when to display the cues (synchronization with the speech signal), where to place the hand (spatial positioning, angle, and size), and how to render it visually (hand design).

The four components are interdependent: timing depends on phoneme alignment, spatial positioning requires both timing and content, and visual rendering builds on all previous stages. This structure is the result of the analysis described above. It defines an architecture for cue generation and supports the implementation of a consistent and extensible system. The same framework has guided the present system and can serve as a reference for future developments.

For example, the system is *multilingual by design* in the sense that language-specific knowledge is externalized into modular, open-format resource files. The core components—covering normalization, phonetic transcription, alignment, and cue generation—are implemented in a language-independent way. Language-specific resources, such as dictionaries, acoustic models, and cueing rules, are handled through separate, editable files. This modular architecture follows the same strategy as adopted in SPPAS for text normalization (Bigi, 2014), phonetic transcription (Bigi, 2016), and alignment (Bigi and Meunier, 2018). Its applicability to multiple languages has already been validated in these components (Lancien et al., 2020; Bigi et al., 2021; Pakrashi et al., 2023), and is here extended to the novel task of Cued Speech generation.

Figure 3 presents the full processing pipeline, from user inputs to the final coded video. It illustrates the modular organization of the system and the sequence of required operations. The first stages involve automatic processing of the input transcript, audio, and video using the open-source SPPAS toolkit (Bigi, 2015), including normalization, phonetization, forced-alignment, and face analysis. These annotations are used without manual correction and provide the foundation for reproducible experiments. The subsequent steps implement the proposed framework, computing the sequence of keys, their temporal and spatial properties, and rendering the virtual hand accordingly.

### 2.1 What to Code

The first component of the system determines the sequence of keys to be produced from the phoneme transcription. Each key encodes a consonant–vowel association as a pair of handshape and position. Based on the aligned phoneme sequence, the system infers the structure and associates each segmental unit with a key of type C-, -V, or CV. A deterministic finite automaton (DFA) formalizes

Figure 3: Workflow of the full process: from the user's data to the coded video

all valid transitions and decomposes complex syllables into successive keys.

This component was previously described and evaluated in a dedicated study (Bigi, 2023). On a manually annotated corpus, the predicted sequences aligned closely with those produced by expert coders, with most deviations reflecting individual preferences rather than systemic errors. The DFA-based system was found to be both reliable and sufficient. A web-based text-to-cue converter[4], developed in collaboration with the deaf community, provides public access to this module for educational and training use.

## 2.2 When to Display the Cues

Once the sequence of keys is defined, the next step is to determine their temporal coordination with speech. It is already known that the hand must anticipate the associated phonemes to allow visual decoding. This principle has been consistently supported in the literature (Cornett, 1967; Bratakos et al., 1998; Duchnowski et al., 1998, 2000) and confirmed by French studies (Cathiard et al., 2003; Attina, 2005; Aboutabit, 2007), which highlight the role of anticipation in perception.

Four timing models were implemented: three drawn from previous work, and a fourth developed specifically for this system. The notation introduced in (Attina, 2005) is used throughout. A1 marks the acoustic onset of the key—consonant onset in 'C-' or 'CV' keys, vowel onset in '-V' keys. A3 marks the acoustic offset—vowel end in 'CV' or '-V' keys, consonant end in 'C-' keys. M1

and M2 represent the start and end of the manual transition. The interval A1–A3 corresponds to the acoustic duration of the key, while M1 and M2 are the time points to be predicted by the models.

Model 1 reproduces the configuration described in (Duchnowski et al., 1998), in which the hand appears 100 ms before the phoneme, with no transition phase. This model was implemented for reference purposes but was not included in the experimental protocol, as later studies (Duchnowski et al., 2000) have shown that Model 2 yields better results. **Model 2** introduces a fixed transition of 150 ms, so that the hand reaches its target 100 ms prior to the phoneme onset.

**Model 3** adjusts anticipation values based on the consonant–vowel structure of the key. It is derived from French-language studies (Attina, 2005), which reported consistent variation in cue timing across key types. Transitions are defined as proportions of the A1–A3 interval, assuming an average duration of 400 ms. For 'CV' and 'C-' keys, M1 starts 62% before A1 and M2 occurs 10% after A1. For '-V' keys, M1 starts 46% before A1 and M2 occurs 21% after A1.

**Model 4** was developed specifically for this system. It extends previous models by incorporating finer adjustments derived from coder expertise and by explicitly modeling transitions involving the neutral position, which are absent from earlier systems. The model adapts timing to speech rate and defines transition points as proportions of the A1–A3 interval.

For the first key, corresponding to a transition from the neutral zone to a facial position, M1 oc-

---

[4]https://auto-cuedspeech.org/textcue.html

11

curs 140% before A1 and M2 20% before A1. For the second key, these values are 125% and 15% before A1. For the third, 100% and 10%. For subsequent keys, M1 is set to 90% and M2 to 5% before A1. For the return to neutral, M1 is delayed to 20% after A1, and M2 to 80% after M1.

## 2.3 Where to Place the Hand in the Video?

This component determines the position, angle, and size of the hand relative to the speaker's face for each frame of the video.

The vowel positions were first defined by expert coders on a theoretical face, then formalized using the 68-point facial landmark model given by SPPAS. Each position is computed as a function of facial landmarks. The formulas used for the positions of French Cued Speech were derived in collaboration with expert coders and adapted to ensure consistency across speakers and morphologies. They are summarized in Table 1 and illustrated in Figure 4. No variability was introduced at this stage: for each frame, the fingertip is placed at the target coordinates.

|     | x = | y = |
| --- | --- | --- |
| **n** | $x_8$ | $y_8 + 4 \cdot (y_8 - y_{57})$ |
| **b** | $x_4 + \frac{1}{2} \cdot |x_{36} - x_0|$ | $y_1 - \frac{1}{3} \cdot |y_1 - y_0|$ |
| **c** | $x_8$ | $y_8 - \frac{1}{5} \cdot |y_8 - y_{57}|$ |
| **m** | $x_{48} - \frac{1}{4} \cdot |x_{48} - x_4|$ | $y_{60}$ |
| **s** | $x_0 - \frac{2}{3} \cdot |x_8 - x_0|$ | $y_4 - \frac{1}{2} \cdot |y_4 - y_3|$ |
| **t** | $x_8$ | $y_8 + 1.2 \cdot |y_8 - y_{57}|$ |

Table 1: Estimated positions from facial landmarks



Figure 4: Estimated positions relatively to the landmarks

Hand orientation is also controlled to improve visual realism. Three models were implemented. **Model 0** uses a fixed angle of 60°, serving as a baseline (Figure 5). **Model 1** introduces expert-defined



Figure 5: Hand angle of Model 0 is 60°.

variations by position. Excluding the neutral zone, the average angle is 71.2°, with a standard deviation of 9.3°. **Model 2** uses a data-driven approach: five annotated frames per position were manually selected to estimate average orientations. It yields an average angle of 61.8° and a standard deviation of 12.5°. Detailed values are given in Table 2.

| Position | Model 1 | Model 2 |
| --- | --- | --- |
| n (chest) | 50° | 50° |
| b (cheek bone) | 75° | 62° |
| c (chin) | 67° | 59° |
| m (mouth) | 73° | 56° |
| s (side) | 83° | 83° |
| t (throat) | 58° | 49° |

Table 2: Hand angles (in degrees) for Models 1 and 2.

The hand size is scaled proportionally to face height and remains fixed throughout. Transitions between positions follow a straight-line trajectory at constant speed. Handshape transitions occur at the midpoint of this trajectory, using a three-frame fade between the two handshapes. These simplifications reflect a design choice: only one spatial parameter is introduced at a time for evaluation.

This component of the system then produces a complete 2D hand trajectory of the hand, it's angle and it's size, for each frame of the given video.

## 2.4 How to Represent the Hand in the Video?

The final module of the system handles the visual rendering of the cueing hand, based on the timing and spatial information computed in the previous stages. This component determines how the hand appears in the video and offers several options in terms of style and visual clarity.

Four handsets were integrated into the system. Two are based on photographs: a male hand set ('l'), and a female hand set ('b') shown in Figure 5.

The other two use 2D illustrations: 'd' displays a uniform yellow shape, while 'c' assigns colors to specific keys to reduce visual confusions—key 3 is pink, key 8 is blue, and the neutral hand is white; all others remain yellow. These assignments build on prior work (Duchnowski et al., 1998) indicating that color can help distinguish keys that are visually similar but phonologically distinct.

Figure 6 shows examples of these handsets, along with enhancement filters described below.



Figure 6: Some hands configurations: "l+1", "l+2", "d+3", "d", "c"

Three visual enhancements were implemented to explore whether additional graphic information could improve the visual distinction between similar handshapes. Each one is exclusive and applies to a single rendering at a time. The first one adds a dot at the fingertip target and a line along the index for keys 3 and 8, to improve distinction from keys 4 and 2, similarly to the 'c' handset. The second one draws a line along the back of the hand and a dot at the target point, highlighting orientation. The 3rd one overlays the full 21-point hand sights with connecting lines, as illustrated in Figure 5.

This rendering module supports both realistic and stylized outputs and can be adapted to user needs or preferences.

## 2.5 System Summary

The system covers the full pipeline of automatic CS generation. Starting from a video, an audio signal, and an orthographic transcript, it performs phoneme alignment, transformation into keys, synchronization of each key with the speech signal, analysis of facial landmarks, determination of hand angle, hand size, handshape transitions, spatial transitions between positions, and visual rendering.

The process results in a synchronized and augmented video, where a virtual hand encodes the Cued Speech transcription with precise timing and positioning. All elements—phonetic inference, timing models, spatial computation, and graphical output—are integrated into a reproducible framework.

This combination of coverage and modularity is, to our knowledge, the first of its kind.

This framework is implemented in Python and released under an open-source license. Its graphical user interface and user-friendly installation process allow non-specialists to use it.

## 3 System Evaluation

The system was evaluated through a decoding task with eight deaf participants, all fluent in French Cued Speech and familiar with video-based cueing. The goal was to assess the readability of automatically generated cues and to compare different configuration options. The task consisted in watching short cued videos and writing down what was decoded. Their responses were scored using SCLite, designed for evaluating ASR output. It aligns each decoded transcription with a reference using utterance IDs and computes word-level scores: correct (Corr), substituted (Sub), deleted (Del), and inserted (Ins). In this setting, the reference is the recorded sentence, and the hypothesis is the participant's transcription.

Decoding accuracy was then used as a proxy for system performance. This metric was deliberately chosen to reflect the perceptual clarity of the generated cues, independently of participant-specific inference or language comprehension skills. Although comprehension-based tasks might better reflect communicative effectiveness, they would confound system output quality with individual-level interpretation strategies. By focusing on transcription alignment, the evaluation isolates the contribution of the system itself, ensuring a more rigorous and interpretable measure of cue readability.

### 3.1 Experimental Conditions

The evaluation was conducted during the 2024 annual internship organized by the ALPC. Eight deaf adults participated on a voluntary basis and gave informed consent. All participants watched a standardized instructional video before the session. The protocol was anonymous, non-intrusive, and approved by the organizing institution.

Each participant decoded 20 silent videos: five manually coded by a professional (used as a reference set), and fifteen automatically generated using the system with different configurations. To control for inter-participant variability, each participant was assigned to a single experimental dimension: timing, angle, hand appearance, or visual enhance-

ment. This allowed for within-subject comparisons across three variants per parameter. Each system configuration was identified by a four-character code: the first digit refers to the timing model (2, 3, or 4), the second to the angle model (0, 1, or 2), the third to hand appearance ('b', 'c', or 'd'), and the fourth to optional enhancements (1, 2 or 3). Participants were divided into four groups:

- **Group 1 – Timing:** P1 and P5 decoded sets 2.1.l.0, 3.1.l.0, and 4.1.l.0.

- **Group 2 – Angle:** P2 and P6 decoded sets 4.0.l.0, 4.1.l.0, and 4.2.l.0.

- **Group 3 – Appearance:** P3 and P7 decoded sets 4.1.b.0, 4.1.c.0, and 4.1.d.0.

- **Group 4 – Enhancement:** P4 and P8 decoded sets 4.1.l.1, 4.1.l.2, and 4.1.l.3.

The five manually coded reference videos were presented first. The fifteen system-generated clips followed, in a fixed interleaved order balancing topic and condition. Playback issues affected two participants (three clips for P1, two for P2) due to local hardware errors. Since all videos had been generated beforehand, only playback was affected and the evaluation protocol remained valid. This is reported here in accordance with FAIR principles.

## 3.2 Global Decoding Performance

Table 3 presents the decoding scores for the control set (professionally coded) and for the system-generated output (all configurations combined). Manual coding achieved 92.3% accuracy. The system, with no participant training or adaptation, reached 80.7%.

| SPK | Corr | Sub | Del | Ins | Err |
|---|---|---|---|---|---|
| Control | 92.3 | 5.2 | 2.5 | 2.3 | 10.0 |
| All sets | 80.7 | 9.7 | 9.6 | 1.3 | 20.6 |

Table 3: Participant decoding scores

These results were obtained using strict word-level alignment. Minor spelling differences were counted as substitutions, and no correction was applied to participant input. The control score reflects the best achievable performance under these conditions and serves as an oracle reference.

That the system reaches over 80% under the same constraints is a key finding. Participants had no prior exposure to the system and received

no training. Despite this, several decoded videos scored near the reference level. The output is therefore not only intelligible but already close to expert quality for a majority of sentences.

The most frequent errors were deletions, increasing from 2.5% in the control set to 9.6% with system output. Substitutions also rose, though to a lesser extent. Informal debriefings suggest that fast speech segments were harder to decode, especially when hand transitions compressed timing contrasts.

To our knowledge, this is the first publicly documented benchmark comparing professional and system-generated Cued Speech. These results show that automatic cue generation is not only feasible, but already yields intelligible output close to expert performance. This first benchmark sets a high baseline for future systems and provides a reproducible framework for comparison.

The 80.7% score reported above reflects an average across multiple system variants. It includes different timing strategies, spatial models, hand appearances, and visual enhancements. The result therefore combines heterogeneous outputs, some of which led to higher decoding scores than others.

## 3.3 Detailed evaluation and discussion

The three sentence sets used in the experiment yielded average decoding scores of 83.6%, 84.4%, and 74.9%, respectively, indicating noticeable differences in difficulty. Without normalization, such variation would interfere with the analysis of model-specific effects. To control for these biases, all scores were normalized by participant and by sentence set. This adjustment accounts for individual decoding ability and for intrinsic difficulty of the material. Final results are reported as $z$-scores: a positive value indicates that the participant decoded better than their own average, and a negative value indicates below-average decoding accuracy.

### 3.3.1 Group 1 – Timing Models

Participant P1 showed slightly negative performance on the baseline (model 2), and slightly positive scores on models 3 and 4 ($z = -0.07$, $+0.04$, $+0.06$). P5 had the best result on model 2 ($z = +0.08$), followed closely by model 4 ($+0.02$), with model 3 performing lower ($-0.04$). Overall, model 4 seems less sensitive to speaker or material, while model 3 is more sensitive to speaker or sentence variation.

### 3.3.2 Group 2 – Angle Models

For P2, model 1 yielded the best performance (+0.03), followed by model 2 (−0.04), while model 0 performed neutrally (−0.001). P6 achieved highest scores on models 0 and 1 (+0.07 and +0.06), with lower performance on model 2 (−0.03). The results suggest that moderate expert-defined angle variation (model 1) provides a good compromise between visual consistency and realism, while corpus-derived angles (model 2) may introduce instability.

### 3.3.3 Group 3 – Hand Appearance

P3 had a slight preference for the 'd' design (+0.01), with lower results on the 'b' and 'c' designs (−0.09, −0.03). P7 also favored 'd' (+0.12), followed by 'b' (+0.05), and had a neutral response to 'c' (−0.01). Unlike earlier findings reported in (Duchnowski et al., 1998), our results do not replicate a consistent benefit from color coding: one participant improved with the 'c' design, while another performed better without it. These observed trends confirm that the simplified, high-contrast 'd' illustrations enhance decoding performance, likely due to their visual clarity and reduced ambiguity.

### 3.3.4 Group 4 – Visual Enhancements

P4 showed balanced performance across the three enhancement types ($z$-scores ranging from 0.0 to +0.04), while P8 experienced a sharp decline, particularly on Skeleton (−0.19). These results suggest that while visual enhancements may assist some users, they may also introduce distracting or overly complex visual elements, especially for less experienced decoders.

### 3.3.5 Discussion

The experimental results converge on a configuration that *favors clarity over realism*. The most effective combination includes a fixed anticipation model refined by phonetic context (Model 4), expert-defined orientation values (Model 1), and a stylized 2D design with strong visual contrast ('d'). This setup does not aim to reproduce natural hand movement but rather to enhance cue discriminability. It consistently produced the best decoding scores across participants and conditions. Visual enhancements overlays did not improve performance and occasionally introduced confusion, suggesting that additional graphic elements may interfere with the perception of essential features. These findings support the adoption of a simpli-fied, controlled rendering strategy as the system's default configuration for future use.

These results highlight that controlled visual simplicity can effectively outperform realism by enhancing usability and reducing cognitive load in accessibility-focused systems.

## 4 Conclusion

Despite the documented benefits of Cued Speech for phonological awareness and literacy, no operational system has yet addressed its automatic generation in a reproducible and open manner. The only prior effort explicitly targeting cue generation in video, developed at MIT in the late 1990s, remains undocumented, non-reproducible, and is no longer maintained.

This paper presents the first functional and publicly available system for automatic Cued Speech generation. It targets French and implements a modular pipeline structured into four components: determining what to code, when to display, where to place, and how to render. Each step is formally defined and operational, from phoneme alignment to video rendering with an integrated virtual hand. The system provides explicit control over linguistic content, synchronization, spatial placement, and visual output.

Evaluation with deaf participants confirmed that the output is readable and effective: decoding accuracy averaged 80.7%, compared to 92.3% for professionally coded videos. This result was obtained without participant training or adaptation. Among the tested parameters, hand appearance had the strongest impact. The highest scores were obtained with a stylized 2D design, limited angle variation, and no visual enhancement. These findings indicate that intelligibility benefits from simplification rather than natural imitation.

This work defines a complete and reproducible framework for Cued Speech generation in video. Moreover, it provides a usable tool with a graphical interface, ready for practical use and offering a reference baseline for future systems. The system is already integrated into the actively maintained software platform SPPAS, and has been successfully used by non-technical users in applied settings.

The next step will involve inserting transitional frames when needed, to reduce deletion errors and improve comfort. The goal is to better match the rhythm of the speaker with the decoding strategies used by human coders.

## Limitations

This study presents the first fully documented and reproducible system for automatic CS generation. However, several limitations must be acknowledged.

First, the system has been implemented and evaluated only for French. While the architecture is designed to support multiple languages, further work is needed to confirm its adaptability to different phonological inventories and cueing conventions. This is currently being addressed through the ongoing adaptation of the system to American English.

Second, although the evaluation protocol was carefully designed, the number of participants remains limited. This constraint, inherent to the difficulty of recruiting expert Cued Speech users, may affect the generalizability of some findings.

Third, while the current design provides transparency and control, it may miss fine-grained variations observed in natural cueing. To address this, a follow-up project has been launched to explore targeted data-driven modeling, restricted to cases where statistical learning is justified — in line with principles of ecological minimalism and methodological necessity.

Finally, two aspects of the system have been fixed *a priori* and remain to be systematically evaluated: the precise spatial placement of hand positions around the face, and the trajectory modeling, which currently assumes straight-line motion at constant speed. While hand cue positions are algorithmically inferred from facial landmarks, we acknowledge that systematic validation against manual annotations remains limited due to the complexity of recruiting trained evaluators. Preliminary cross-checks on held-out data indicate promising consistency, and ongoing work is extending this evaluation as resources permit.

## Ethical Considerations

This study did not involve the collection of any sensitive or identifying information. Participation was voluntary, based on informed consent, and fully anonymous. Participants were not evaluated; rather, their responses served solely to assess the intelligibility of the system's outputs.

The experiment followed the principles of the ALPC association's internal ethics charter, which promotes respect, autonomy, and non-discrimination in all interactions with deaf participants and their families.

## Acknowledgements

## References

N. Aboutabit. 2007. *Reconnaissance de la Langue Française Parlée Complété (LPC): décodage phonétique des gestes main-lèvres.* Ph.D. thesis, Institut National Polytechnique de Grenoble - INPG.

V. Attina. 2005. *La Langue Française Parlée Complétée: Production et Perception.* Ph.D. thesis, Institut National Polytechnique de Grenoble - INPG.

B. Bigi. 2014. A multilingual text normalization approach. *Human Language Technology Challenges for Computer Science and Linguistics, LNAI 8387*, pages 515–526.

B. Bigi. 2015. SPPAS - Multi-lingual Approaches to the Automatic Annotation of Speech. *The Phonetician*, 111–112:54–69.

B. Bigi. 2016. A phonetization approach for the forced-alignment task in SPPAS. *Human Language Technology. Challenges for Computer Science and Linguistics, LNAI 9561*, pages 515–526.

B. Bigi. 2023. An analysis of produced versus predicted french cued speech keys. In *10th Language & Technology Conference: Human Language Technologies as a Challenge for Computer Science and Linguistics, ISBN: 978-83-232-4176-8*, pages 24–28, Poznań, Poland.

B. Bigi and C. Meunier. 2018. Automatic segmentation of spontaneous speech. *Revista de Estudos da Linguagem. International Thematic Issue: Speech Segmentation*, 26(4).

B. Bigi, A-S. Oyelere, and B. Caron. 2021. Resources for automated speech segmentation of the african language naija (nigerian pidgin). *Human Language Technology. Challenges for Computer Science and Linguistics, LNAI 12598*, pages 164–173.

B. Bigi, M. Zimmermann, and C. André. 2022. CLeLfPC: a Large Open Multi-Speaker Corpus of French Cued Speech. In *Proceedings of The 13th Language Resources and Evaluation Conference*, pages 987–994, Marseille, France.

---

M. S. Bratakos. 1995. *The effect of imperfect cues on the reception of cued speech*. Ph.D. thesis, Massachusetts Institute of Technology.

M. S. Bratakos, P. Duchnowski, and L. D. Braida. 1998. Toward the automatic generation of cued speech. *Cued Speech Journal*, 6:1–37.

M.-A. Cathiard, V. Attina, and D. Alloatti. 2003. Labial anticipation behavior during speech with and without cued speech. In *Proceedings of the 15th International Congress of Phonetic Sciences*, pages 1935–1938, Barcelona, Spain.

B. R. Clarke and D. Ling. 1976. The effects of using cued speech: A follow-up study. *Volta Review*, 78(1):23–34.

S. Colin, J. Ecalle, E. Truy, G. Lina-Granade, and A. Magnan. 2017. Effect of age at cochlear implantation and at exposure to cued speech on literacy skills in deaf children. *Research in developmental disabilities*, 71:61–69.

R. O. Cornett. 1967. Cued speech. *American annals of the deaf*, pages 3–13.

R. O. Cornett, R. Beadles, and B. Wilson. 1977. Automatic cued speech. In *Research Conference on Speech Processing Aids for the Deaf*, pages 224–239, Gallaudet College (USA).

P. Duchnowski, L. D. Braida, D. Lum, M. Sexton, J. Krause, and S. Banthia. 1998. Automatic generation of cued speech for the deaf: status and outlook. In *International Conference on Auditory-Visual Speech Processing*, Sydney, Australia.

P. Duchnowski, D. S. Lum, J. C. Krause, M. G. Sexton, M. S. Bratakos, and L. D. Braida. 2000. Development of speechreading supplements based on automatic speech recognition. *IEEE transactions on biomedical engineering*, 47(4):487–496.

C. G Fisher. 1968. Confusions among visually perceived consonants. *Journal of speech and hearing research*, 11(4):796–804.

M. Lancien, M-H. Côté, and B. Bigi. 2020. Developing resources for automated speech processing of quebec french. In *Proceedings of The 12th Language Resources and Evaluation Conference*, pages 5323–5328, Marseille, France. European Language Resources Association.

C. J. LaSasso, K. L. Crain, and J. Leybaert. 2010. *Cued speech and cued language development for deaf and hard of hearing children*. Plural Publishing.

J. Leybaert and B. Charlier. 1996. Visual speech in the head: The effect of cued-speech on rhyming, remembering, and spelling. *The Journal of Deaf Studies and Deaf Education*, 1(4):234–248.

D. W. Massaro and Stephen E. Palmer Jr. 1998. *Perceiving talking faces: From speech perception to a behavioral principle*. Mit Press.

N. A Neef and B. A. Iwata. 1985. The development of generative lipreading skills in deaf persons using cued speech training. *Analysis and intervention in developmental disabilities*, 5(4):289–305.

M. Pakrashi, B. Bigi, and S. Mahanta. 2023. Resources creation of bengali for sppas. In *10th Language & Technology Conference: Human Language Technologies as a Challenge for Computer Science and Linguistics, ISBN: 978-83-232-4176-8*, pages 218–222, Poznań, Poland.

J. Rönnberg. 1995. Perceptual compensation in the deaf and blind: Myth or reality? *Compensating for psychological deficits and declines: Managing losses and promoting gains*, pages 251–274.

J. Rönnberg, S. Samuelsson, B. Lyxell, R Campbell, B Dodd, and D Burnham. 1998. Conceptual constraints in sentence-based lipreading in the hearing-impaired. *The psychology of speechreading and auditory–visual speech*, pages 143–153.

M. G. Sexton. 1997. *A video display system for an automatic cue generator*. Ph.D. thesis, Massachusetts Institute of Technology.

## A Reproducibility

All data and source code referenced in this paper comply with the principles of open science. The source code of the proposed system is released under the GNU Affero General Public License v3 (AGPLv3). It is part of SPPAS and can be downloaded at `https://sourceforge.net/projects/sppas/`.

The experimental scripts are also made available under the same license and can be obtained from the author upon request.

The datasets used in this work are distributed under both the Open Database License v1.0 (ODbL) and the Creative Commons Attribution-NonCommercial 4.0 International (CC BY-NC 4.0) licenses. They can be downloaded at `https://hdl.handle.net/11403/clelfpc/v10`.

**Software and Evaluation Tools:**

- The full speech segmentation pipeline, including text normalization, phonetic transcription, and alignment, was performed using SPPAS, version 4.11 (`https://sppas.org/`),

- Evaluation metrics were computed using SCTK 2.4.12 (`https://github.com/usnistgov/SCTK`).

# `topicwizard` - a Modern, Model-agnostic Framework for Topic Model Visualization and Interpretation

**Márton Kardos**
Aarhus University
martonkardos@cas.au.dk

**Kenneth C. Enevoldsen**
Aarhus University
kenneth.enevoldsen@cas.au.dk

**Kristoffer Laigaard Nielbo**
Aarhus University
kln@cas.au.dk

## Abstract

Topic models are statistical tools that allow their users to gain qualitative and quantitative insights into the contents of textual corpora without the need for close reading (Nielbo et al., 2024). They can be applied in a wide range of settings from discourse analysis (Bednarek, 2024), through pretraining data curation (Peng et al., 2025), to text filtering (Ma et al., 2016). Topic models are typically parameter-rich, complex models, and interpreting these parameters can be challenging for their users. It is typical practice for users to interpret topics based on the top 10 highest ranking terms on a given topic. This *list-of-words* approach, however, gives users a limited and biased picture of the content of topics (Gillings and Hardie, 2022). Thoughtful user interface design and visualizations can help users gain a more complete and accurate understanding of topic models' output. While some visualization utilities do exist for topic models, these are typically limited to a certain type of topic model. We introduce `topicwizard` [1], a framework for model-agnostic topic model interpretation, that provides intuitive and interactive tools that help users examine the complex semantic relations between documents, words and topics learned by topic models.

## 1 Introduction

Topic models are statistical instruments, which have been developed to uncover human-interpretable topics in corpora of text (Blei, 2012). These methods have allowed analysts gain insights into the contents of large corpora, the manual reading of which would be impractical or impossible. Topic models also often offer a more impartial account of a corpus' content (Nielbo et al., 2024).

Typically, topic models' outputs are presented to users in the form of the highest-ranking words and

perhaps documents on a given topic. While this allows users to gain a superficial understanding of a topic, one might miss crucial details, and a lot of nuances, when topic models are exmined this way (Gillings and Hardie, 2022). We suggest that topic models capture more detailed information about topics than simple word lists convey, and that carefully designed interfaces can help users better explore this complexity.

### 1.1 Topic Models are Diverse

While topic models all carry out a similar task, they can also be very different from each other in how they conceptualize topic discovery.

Topic models originally relied on a bag-of-words model of documents where they are represented as sparse vectors of word-occurrence counts, with an optionally applied weighting scheme, such as tf-idf. Most commonly, these models either discover topics by matrix factorization (Gillis and Vavasis 2014, Kherwa, Pooja and Bansal, Poonam 2017) or by fitting a probabilistic generative model over these representations (Blei et al. 2003, Yin and Wang 2014, Hofmann 1999) or biterms (Yan et al., 2013).



Figure 1: A Simplified Taxonomy of Topic Models

More recent topic models, however, also rely on context-sensitive, dense text representations from neural networks (Reimers and Gurevych, 2019). These models can conceptualize topic discovery as document clustering and post-hoc term importance

---

[1] https://github.com/x-tabdeveloping/topicwizard

estimation (Grootendorst 2022, Angelov 2020), document generation with amortized variational inference (autoencoders) (Bianchi et al. 2021a, Bianchi et al. 2021b), semantic relation reconstruction (Wu et al., 2024), or semantic decomposition (Kardos et al. 2025a, Kristensen-McLachlan et al. 2024).

## 1.2 Topic Models are Alike

Despite these differences, all topic models have a lot in common. Topic models, in essence, learn a three-way relationship between `words`, `documents` and `topics`.



Figure 2: Common Components Computed by Topic Models

All topic models have a method for extracting the $K$ most relevant words from the discovered topics. These top $K$ words are calculated from a **topic-term matrix** ($\phi$), which is either inferred as part of topic discovery. This matrix has $N$ rows, corresponding to the number of topics, and $M$ columns corresponding to the size of the model's vocabulary. In addition, models compute a **document-topic-matrix** ($\Theta$), where rows represent the $D$ documents in the corpus, while the $N$ columns represent topics. This matrix contains the importance/relevance of a topic in a document.

## 1.3 Contribution

We introduce `topicwizard`, a model-agnostic topic model visualization framework that allows users to investigate complex semantic relations between words, documents and topics in their corpora. `topicwizard` is natively compatible with topic modelling libraries, which use the scikit-learn API (Pedregosa et al., 2011), such as tweetopic (Kardos, 2022) and Turftopic (Kardos et al., 2025b) and comes with compatibility layers for Gensim and BERTopic.

## 2 Related Work

Due to Latent Dirichlet Allocation's (LDA) popularity, a considerable amount of work has been dedicated to visualizing and interpreting its outputs. Chuang et al. (2012b) discuss best practices and design considerations for visualization and interpretation systems for LDA. Chuang et al. (2012a) introduced the Termite system for interactively visualizing and interpreting LDA output. The main visualization in Termite is a stylized version of the topic-term matrix (see Figure 8), where circles of different size are at the intersection of terms and topics indicating their importance. The authors also propose a scheme for selecting the most topically salient words, since displaying all words in the corpus would not be feasible. As a consequence, Termite can only display a limited number of words. Additionally, Termite is no longer under active maintenance [2].

LDAvis (Sievert and Shirley, 2014) is an interactive visualization R package for LDA (see Figure 9). LDAvis combines elements of previous topic visualization systems, including an inter-topic distance map, term distribution plots, and a term-weighting scheme to show only the most specific and (*relevant*) terms. Similar to Termite, the original LDAvis package is no longer maintained. Its Python port, PyLDAvis, receives occasional updates, but does not enjoy feature parity with the original package.

Notable visualization utilities are also included in the BERTopic library (Grootendorst, 2022), which boasts model-specific plotting functions, such as an inter-topic map, document cluster visualizations, and term distribution bar-charts. Similarly, Turftopic (Kardos et al., 2025b) also contains model-specific visualization utilites for a number of models, including cluster maps, concept compasses for $S^3$ (Kardos et al., 2025a) and interactive timeline plots for dynamic topic models. While these visualizations are useful, they are typically of limited interactivity, and are limited to a particular type of model.

## 3 🔲 `topicwizard`

To address these challenges, we outline `topicwizard`, a novel system for topic model interpretation. Our framework is model-agnostic,

---

[2]The Termite repository on Github was last committed to 11 years prior to the writing of this article

Figure 3: An overview of visualizations and pages in the `topicwizard` framework
*All visualizations were produced using KeyNMF ([Kristensen-McLachlan et al., 2024](#))*

allows users to investigate topic models from a number of distinct perspectives, and is highly interactive, thereby providing a more complete picture of topic models' output,

## 3.1 Topic Models Learn Topic Representations

Topic models' primary objective is to discover latent themes in a corpus. Being able to understand what concepts make up such topics, and how these topics are related is perhaps the most important aspect of interpreting topic models.

In `topicwizard` (see Figure 3a), similar to [Sievert and Shirley (2014)](#) an `inter-topic map` is displayed, which shows the relative distances of topics to each other. While [Sievert and Shirley (2014)](#) utilize PCA for this visualization, projections in `topicwizard` are calculated with UMAP ([McInnes et al., 2018](#)), since it is better at capturing local structure. The size of the topics on the graph is determined by a *topic importance* score. This score, and thereby the size on the graph indicates how prevalent a given topic is in the corpus overall, also taking into account the length of the documents. Topic importance is calculated in the following manner:

$$s_t = \sum_d^D \Theta_{dt} \cdot |d|$$

where $\Theta_{dt}$ is the importance of topic $t$ and document $d$ and $|d|$ is the number of terms in a given document, and $D$ is the size of the corpus.

To provide users with insights about topics' word content, the `topic-word plot` displays the distribution of the highest ranking words for a given topic, and also how globally prevelant these words are across topics [3]. Since 10-20 words are rarely enough to give a complete picture of the words relevant to a topic, a more comprehensive `topic wordcloud` is also displayed To aid further analysis, users can also manually name topics on this page.

## 3.2 Topic Models Learn Word Embeddings

While topic models' are mainly oriented at discovering topics, they also implicitly learn meaningful representation of words within the corpus. Each column of the topic-term matrix can technically be thought of as a semantic embedding for a given word, with the dimensions being interpretable. This implicit learning of word representations allows us to examine words' relation to each other in a corpus, without explicit reference to the topics.

In `topicwizard` (see Figure 3c), a `word map` is displayed to users, allowing them to quickly and interactively investigate the semantic landscape of words in their corpus. Word positions are calculated by projecting word embeddings to two dimensions using UMAP.

Word embeddings are useful for investigating associative relations in corpora, and have been used for a variety purposes such as query expansion

---

[3]Unlike LDAvis, we do not compute *relevance* scores, since they rely on the assumption that $\phi$ contains word probabilities.

(Kuzi et al., 2016), or to uncover authorship patterns in literature (Baunvig, 2024). Clicking on a word on the word map highlights the words most closely related to the selected one and displays the topical distribution of the selected term and its neighbourhood on the `word-topic plot`. Displaying closely associated words with the selected keywords in topic models can give practitioners a more nuanced picture of word use (Liu and Lei, 2018).

### 3.3 Topic Models Organize Documents

An important aspect of topic models is that they learn a representation of documents in the corpus they are fitted on. Document representations discovered by topic models were historically used for a number of purposes, including retrieval (Yi and Allan, 2009), and studying information dynamics (Barron et al., 2018).

In `topicwizard` (see Figure 3d), a `document map` is displayed, where document's UMAP-projected embeddings can be seen, and documents are coloured based on most prevalent topic. In the case of BoW models, these representations are derived from the document-topic matrix, while with contextual models, the pre-computed sentence embeddings are used.

Secondly, individual documents' contents can be investigated on a `document-topic plot`, which displays the distribution of the most relevant topics, a `document-topic timeline`, which displays how the topical content changes throughout the course of the document and a `document viewer`, where a snippet of the document is displayed, and the most topically relevant words are highlighted. The combination of these document inspection utilities can help users ground and verify topic models' output in the documents themselves, which elevates trust (Chuang et al., 2012b). Additionally, this interface encourages close reading, which provides additional insight into the corpus' content.

### 3.4 Topics Augment User-Defined Groups

Commonly, users of topic models also have some externally defined grouping of documents, which might be relevant for their analyses. This could be binning documents by time period, predefined categories or place of origin. While most topic models do not utilize external labels, meaningful inferences can be made about topics' relation to these labels post-hoc.

An important part of this process is to compute a

**group-topic matrix**, the cells of which contain the summed importance of a given topic for documents in a given group:

$$G_{ij} = \sum_{k}^{D} \Theta_{kj} \cdot I(g_k = i)$$

where $G_{ij}$ is the importance of group $i$ for topic $j$, $g_k$ is the group label of document $k$, and $I(g_k = i)$ is the indicator function.

In `topicwizard` (see Figure 3b), semantic distances between user-defined groups can be seen on the `group map`, where group-topic representations are projected to 2D space using UMAP. Groups are coloured based on the dominant topic in the group. Topic distributions in groups can be seen on the `group-topic plot`, and groups' lexical content can be examined in detail on the `group wordcloud` to the right.

### 3.5 Software Design Considerations

The `topicwizard` Python package was designed with both research and enterprise use in mind. As such, our goal was to develop a package that is accessible to new users and sufficiently flexible to accommodate specific use cases – ranging from academic writing and technical reporting to enabeling business analysts to interact with topic models via a web interface.

The **Web Application** (see Figures 4 - 7) was designed to make topic model interpretation as seemless and quick as possible, in as many environments as possible, including Jupyter notebooks, in the browser, or deployed to the cloud. which produces a readily deployable Docker project to a specified folder.

The **Figures API** makes it trivial for our users to produce specific figures tailored to their needs. This is especially crucial for producing publications, since some colour schemes, fonts or aspect ratios, while appropriate for an interactive web application, might not be visually appealing in a static document.

## 4 Conclusion

This paper introduces `topicwizard`, a comprehensive, interactive, and model-agnostic topic model visualization framework. Our framework is a notable extension over previous topic model visualization systems, thanks to a) supporting a much wider range of models b) allowing users to ground topic models in the corpus, and investigate them from

numerous angles and c) being flexible, actively supported, and production-ready. The `topicwizard` software package has so far been downloaded more than 45000 times from PyPI, demonstrating that practitioners have already found it useful.

## Limitations

While `topicwizard` is the most comprehensive topic model visualization tool to date, it still lacks coverage of a number of aspects of topic modelling. It, for instance, does not have visualization utilities for dynamic, hierarchical and supervised topic models. This is a clear limitation and will have to be addressed in future package releases.

Our framework, as of now, does not provide any utilities for comparing outputs from different topic models either. This is yet another aspect that future work should address.

Furthermore, while we consider model-angosticity to be one of the strengths of our approach, it does, to an extent, limit its usefulness for certain models. Certain visualizations, such as concept compasses, might be highly useful tools for examining the output of Semantic Signal Separation, but their utility might be limited for clustering topic models. We encourage our users, therefore, to use `topicwizard` in tandem with model-specific interpretation utilities from libraries such as BERTopic or Turftopic.

## References

Dimo Angelov. 2020. Top2vec: Distributed representations of topics. *Preprint*, arXiv:2008.09470.

Alexander T. J. Barron, Jenny Huang, Rebecca L. Spang, and Simon DeDeo. 2018. Individuals, institutions, and innovation in the debates of the french revolution. *Proceedings of the National Academy of Sciences*, 115(18):4607–4612.

Katrine Frøkjær Baunvig. 2024. *Grundtvig og Monstrene*. Center for Grundtvigforskning, Aarhus Universitet.

Monika Bednarek. 2024. Topic modelling in corpus-based discourse analysis: Uses and critiques. *Discourse Studies*.

Federico Bianchi, Silvia Terragni, and Dirk Hovy. 2021a. Pre-training is a Hot Topic: Contextualized Document Embeddings Improve Topic Coherence. In *Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing (Volume 2: Short Papers)*, pages 759–766, Online. Association for Computational Linguistics.

Federico Bianchi, Silvia Terragni, Dirk Hovy, Debora Nozza, and Elisabetta Fersini. 2021b. Cross-lingual Contextualized Topic Models with Zero-shot Learning. In *Proceedings of the 16th Conference of the European Chapter of the Association for Computational Linguistics: Main Volume*, pages 1676–1683, Online. Association for Computational Linguistics.

David M. Blei. 2012. Probabilistic topic models. *Commun. ACM*, 55(4):77–84.

David M. Blei, Andrew Y. Ng, and Michael I. Jordan. 2003. Latent dirichlet allocation. *J. Mach. Learn. Res.*, 3:993–1022.

Jason Chuang, Christopher D. Manning, and Jeffrey Heer. 2012a. Termite: visualization techniques for assessing textual topic models. In *Proceedings of the International Working Conference on Advanced Visual Interfaces*, AVI '12, page 74–77, New York, NY, USA. Association for Computing Machinery.

Jason Chuang, Daniel Ramage, Christopher Manning, and Jeffrey Heer. 2012b. Interpretation and trust: designing model-driven visualizations for text analysis. In *Proceedings of the SIGCHI Conference on Human Factors in Computing Systems*, CHI '12, page 443–452, New York, NY, USA. Association for Computing Machinery.

Mathew Gillings and Andrew Hardie. 2022. The interpretation of topic models for scholarly analysis: An evaluation and critique of current practice. *Digital Scholarship in the Humanities*, 38(2):530–543.

Nicolas Gillis and Stephen A. Vavasis. 2014. Fast and Robust Recursive Algorithmsfor Separable Nonnegative Matrix Factorization. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 36(4):698–714.

Maarten Grootendorst. 2022. Bertopic: Neural topic modeling with a class-based tf-idf procedure. *Preprint*, arXiv:2203.05794.

Thomas Hofmann. 1999. Probabilistic latent semantic indexing. In *Proceedings of the 22nd Annual International ACM SIGIR Conference on Research and Development in Information Retrieval*, SIGIR '99, page 50–57, New York, NY, USA. Association for Computing Machinery.

Márton Kardos, Jan Kostkan, Kenneth Enevoldsen, Arnault-Quentin Vermillet, Kristoffer Nielbo, and Roberta Rocca. 2025a. $S^3$ - Semantic Signal Separation. In *Proceedings of the 63rd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 633–666, Vienna, Austria. Association for Computational Linguistics.

Márton Kardos. 2022. tweetopic: Blazing fast topic modelling for short texts. *GitHub repository*.

Márton Kardos, Kenneth C. Enevoldsen, Jan Kostkan, Ross Deans Kristensen-McLachlan, and Roberta

Rocca. 2025b. Turftopic: Topic modelling with contextual representations from sentence transformers. *Journal of Open Source Software*, 10(111):8183.

Kherwa, Pooja and Bansal, Poonam. 2017. Latent semantic analysis: An approach to understand semantic of text. In *2017 International Conference on Current Trends in Computer, Electrical, Electronics and Communication (CTCEEC)*, pages 870–874.

Ross Deans Kristensen-McLachlan, Rebecca Marie Matouschek Hicke, Márton Kardos, and Mette Thunø. 2024. Context is Key(NMF): Modelling Topical Information Dynamics in Chinese Diaspora Media. In *Proceedings of the Computational Humanities Research Conference 2024*, volume 3834 of *CEUR Workshop Proceedings*, pages 829–847, Germany. CEUR-WS.

Saar Kuzi, Anna Shtok, and Oren Kurland. 2016. Query expansion using word embeddings. In *Proceedings of the 25th ACM International on Conference on Information and Knowledge Management*, CIKM '16, page 1929–1932, New York, NY, USA. Association for Computing Machinery.

Dilin Liu and Lei Lei. 2018. The appeal to political sentiment: An analysis of donald trump's and hillary clinton's speech themes and discourse strategies in the 2016 us presidential election. *Discourse, Context & Media*, 25:143–152.

Jialin Ma, Yongjun Zhang, Jinling Liu, Kun Yu, and XuAn Wang. 2016. Intelligent sms spam filtering using topic model. In *2016 International Conference on Intelligent Networking and Collaborative Systems (INCoS)*, pages 380–383.

Leland McInnes, John Healy, Nathaniel Saul, and Lukas Großberger. 2018. Umap: Uniform manifold approximation and projection. *J. Open Source Softw.*, 3:861.

Kristoffer L. Nielbo, Folgert Karsdorp, Melvin Wevers, Alie Lassche, Rebekah B. Baglini, Mike Kestemont, and Nina Tahmasebi. 2024. Quantitative text analysis. *Nature Reviews Methods Primers*, 4(1).

F. Pedregosa, G. Varoquaux, A. Gramfort, V. Michel, B. Thirion, O. Grisel, M. Blondel, P. Prettenhofer, R. Weiss, V. Dubourg, J. Vanderplas, A. Passos, D. Cournapeau, M. Brucher, M. Perrot, and E. Duchesnay. 2011. Scikit-learn: Machine learning in Python. *Journal of Machine Learning Research*, 12:2825–2830.

Jiahui Peng, Xinlin Zhuang, Qiu Jiantao, Ren Ma, Jing Yu, Tianyi Bai, and Conghui He. 2025. Unsupervised Topic Models are Data Mixers for Pre-training Language Models. *Preprint*, arXiv:2502.16802.

Nils Reimers and Iryna Gurevych. 2019. Sentence-BERT: Sentence Embeddings using Siamese BERT-Networks. In *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)*, pages

3982–3992, Hong Kong, China. Association for Computational Linguistics.

Carson Sievert and Kenneth Shirley. 2014. LDAvis: A method for visualizing and interpreting topics. In *Proceedings of the Workshop on Interactive Language Learning, Visualization, and Interfaces*, pages 63–70, Baltimore, Maryland, USA. Association for Computational Linguistics.

Xiaobao Wu, Thong Thanh Nguyen, Delvin Ce Zhang, William Yang Wang, and Anh Tuan Luu. 2024. FASTopic: Pretrained Transformer is a Fast, Adaptive, Stable, and Transferable Topic Model. In *The Thirty-eighth Annual Conference on Neural Information Processing Systems*.

Xiaohui Yan, Jiafeng Guo, Yanyan Lan, and Xueqi Cheng. 2013. A biterm topic model for short texts. In *Proceedings of the 22nd International Conference on World Wide Web*, WWW '13, page 1445–1456, New York, NY, USA. Association for Computing Machinery.

Xing Yi and James Allan. 2009. A comparative study of utilizing topic models for information retrieval. In *Advances in Information Retrieval*, pages 29–41, Berlin, Heidelberg. Springer Berlin Heidelberg.

Jianhua Yin and Jianyong Wang. 2014. A dirichlet multinomial mixture model-based approach for short text clustering. In *Proceedings of the 20th ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*, KDD '14, page 233–242, New York, NY, USA. Association for Computing Machinery.

# A  Appendix

See Figures 4-7 for screenshots of `topicwizard`, Figure 9 for LDAvis and Figure 8 for Termite.



Figure 4: Screenshot of the `Topics` page in the `topicwizard` Web Application



Figure 5: Screenshot of the `Words` page in the `topicwizard` Web Application

Figure 6: Screenshot of the `Documents` page in the `topicwizard` Web Application



Figure 7: Screenshot of the `Groups` page in the `topicwizard` Web Application

Figure 8: Screenshot of the Termite System
Figure from (Chuang et al., 2012a)



Figure 9: Screenshot of LDAvis
Figure from (Sievert and Shirley, 2014)

# Scalable Text Vectorization with Hyperdimensional Computing Through Selective Word Encoding

**Timur Mudarisov**
University of Luxembourg
Luxembourg
timur.mudarisov@uni.lu

**Zsofia Kraussl**
Bayes Business School
London
zsofia.kraussl@bayes.ac.uk

**Evgeny Polyachenko**
University of Luxembourg
Luxembourg
evgeny.polyachenko@uni.lu

**Enriqueta Patricia Becerra Sanchez**
University of Luxembourg
Luxembourg
enriqueta.becerra@uni.lu

**Tatiana Petrova**
University of Luxembourg
Luxembourg
tatiana.petrova@uni.lu

**Radu State**
University of Luxembourg
Luxembourg
radu.state@uni.lu

## Abstract

Hyperdimensional Computing (HDC) is a promising approach for various machine learning tasks. In this work, we focus on its application to encoding large text datasets, where the *curse of dimensionality* presents a significant challenge. To mitigate this issue, we employ compression techniques that are based on classical models such as Term Frequency-Inverse Document Frequency (TF-IDF) and Latent Dirichlet Allocation (LDA). We derive theoretical expressions for Compression Rate, Jensen-Shannon Divergence, and ROUGE score, which quantify text size reduction, preservation of word distributions, and retention of key information, respectively. These expressions are validated using the IMDB, arXiv, and AG News datasets. Our results demonstrate that TF-IDF compression can reduce the encoded text size to 10% (or less in some cases) of the original input while also achieving slightly worse distinguishability between classes in classification tasks.

## 1 Introduction

Hyperdimensional Computing (HDC) is a machine learning approach inspired by principles of neural computation. It represents and manipulates data through high-dimensional vectors, typically in the order of thousands or millions, enabling information processing and storage. This methodology exhibits inherent robustness to noise, offers efficient learning capabilities, and effectively handles complex, unstructured data (Kanerva, 2009). HDC has gained considerable interest in emerging applications, such as robotics and health diagnostics, alongside established areas including data center recommendation systems (Mitrokhin et al., 2019; Neubert and Schubert, 2021; Yunhui et al., 2021). This increasing adoption and interest highlights the need for a robust theoretical justification. To address this, researchers have investigated HDC from different perspectives. These studies include an in-depth examination of its geometric characteristics (Pourmand et al., 2024), a comprehensive analysis of its algebraic foundations (Yu et al., 2024), and a detailed investigation of encoding structures used within HDC systems (Thomas et al., 2021). Each perspective contributes to a deeper understanding of HDC and its potential applications.

Kanerva (2009) identified several valuable aspects of different HDC realizations. These include their robustness to noise, which allows HDC to maintain performance despite disruptions. Their inherent transparency also helps the understanding and interpretation of results. Furthermore, HDC exhibit useful distributed properties, which enable efficient parallel processing, for example using GPUs. HDC have been successfully applied in various scientific fields (Rahimi et al., 2019; Kanerva, 2009), and their application to Natural Language Processing (NLP) tasks is of particular interest. Specifically, Kleyko et al. (2023) demonstrated successful applications of HDC to translation, sentence similarity, and topic classification problems. However, Thomas et al. (2021) pointed out important limitations of basic HDC. Among these, a critical challenge is the *curse of dimensionality*. This effect describes how increases in data size can cause an exponential rise in vector space dimensionality, complicating analysis and processing.

To address the challenge of the curse of dimen-

sionality in HDC, we propose using text compression techniques. In this paper, we aim at exploring two classical techniques for text compression: TF-IDF selection (Spärck Jones, 1972) and LDA (Blei et al., 2003). Our contribution to the state-of-the-art in HDC is threefold: First, in Section 2 we introduce a novel model – *compression HDC* (CHDC) which combines a theory-based encoding procedure with data compression using TF-IDF or LDA. This model allows encoding information efficiently while reducing the size of representations. Second, we analyze the compression effect of these techniques (Section 3.1), providing theorems that estimate the compression rate. Third, we examine the encoding effect of the binary uniform HDC (Section 3.2), showing that our results are robust to different conditions. In Section 4, we experimentally validate our theoretical findings, for the quality of the proposed compression and encoding processes. Finally, Section 5 wraps up and discusses prospects.

## 2 Model Setup

The scheme of our proposed model is presented in Figure 1. Before any text analysis is performed, a standard procedure of pre-processing is used and is therefore not shown in the scheme. This procedure involves four steps applied to a large text (document): first, only letters and numbers are retained; second, the text is broken down into words; third, lemmatization is applied, which reduces words to their base or dictionary form (lemma); and finally, stemming is applied, which reduces words to their root form.



Figure 1: Workflow of the compression HDC model, illustrating the processing of a large text using text compression and HDC encoding (blue), to produce a final embedding.

The core of our proposed compression HDC model is defined by two components: compression and HDC encoding. These components are detailed in Sections 2.1 and 2.2, respectively.

### 2.1 Compression procedure

Let $\mathcal{W} = \{w_1, \ldots, w_M\}$ represent a set of $M$ unique words and corpus $\mathcal{D} = \{d_1, \ldots, d_N\}$ is a set of $N$ documents. Given these sets $(\mathcal{W}, \mathcal{D})$, our

goal is to reduce the number of words in each document by focusing on the most informative ones. To achieve this, we assign a score to each word and extract the set of word-score pairs $\{(w, s_w)\}$. For the TF-IDF-based compression, we define the score as follows:

**Definition 1.** *The TF-IDF score for a word $w_i$ in a corpus $\mathcal{D}$ is defined as:*

$$s_w = \mathrm{ts}(w, \mathcal{D}) = \frac{1}{N} \sum_{j=1}^{N} f_{w,j} \ln \frac{N}{N_w}, \qquad (1)$$

*where $f_{w,j}$ is the frequency of word $w$ in document $d_j$, $N_w$ is the number of documents in $\mathcal{D}$ containing word $w$.*

Note that our definition differs from the standard TF-IDF definition, which depends on $w$, $d$ and $\mathcal{D}$ and does not contain averaging over documents.

Latent Dirichlet Allocation (LDA) assumes that documents are represented as bags of words, where each document is a mixture of $T$ topics, with $T$ being a predefined number of topics. The probability of a word $w$ belonging to topic $t$ is denoted as $\phi_{t,w}$. The matrix $\Phi = \{\phi_1, \ldots, \phi_T\} \in \mathbb{R}^{T \times M}$, where each $\phi_t$ represents the probability distribution of words for topic $t$, is determined by maximizing the likelihood function $\mathbb{P}(\mathcal{W}, \mathcal{D}|\Phi, \alpha)$, and $\alpha \in \mathbb{R}_+^T$ are the parameters of the Dirichlet distribution (Blei et al., 2003). Based on the LDA model, we define the score as follows:

**Definition 2.** *The LDA-based score for a word $w$ in topic $t$ is defined as:*

$$s_{t,w} = \phi_{t,w}. \qquad (2)$$

We consider the documents unordered and refer to them interchangeably using either the index $j$ or the document $d$ itself, as an element of $\mathcal{D}$. For words and word-related quantities, we will refer to them interchangeably using either the word $w$ itself or the index $i$, specifying the ordering when necessary. Thus, for example, $f_{i,j}$ and $f_{w,d}$ denote the same quantity.

We present the following *compression criteria*. For TF-IDF-based compression, we select the $p$-quantile of words with the highest scores from the set $\{(w, s)\}_{w \in \mathcal{W}}$, resulting in a reduced dictionary $\mathcal{W}_p$ containing approximately $pM$ words. For LDA-based compression, we select the top $pM$ words from each topic, based on their topic probabilities $s_{w,t}$. Because each word has a probability of belonging to every topic, the resulting

reduced dictionary $\mathcal{W}_p$ typically contains fewer than $TpM$ words. Subsequently, we create a new set of compressed documents $\mathcal{D}' = \{d'_1, \ldots, d'_N\}$, where each $d'_j$ is formed by combining words from $\mathcal{W}_p$, preserving the most important words of the original document and their sequential order within each document.

To evaluate the compression quality, we introduce three classical performance metrics:

1. **Compression rate**. A standard metric in compression theory, defined as the ratio:

$$\text{CR} = \frac{\sum_{j=1}^{N} |d'_j|}{\sum_{j=1}^{N} |d_j|} \, , \qquad (3)$$

   where $|d_j|$ and $|d'_j|$ denote the total number of non-unique words in the uncompressed document $d_j$ and the compressed document $d'_j$, respectively. This metric directly quantifies the reduction in text size achieved by compression.

2. **Jensen-Shannon divergence.** For distributions $p$ and $q$, the Jensen-Shannon divergence (JSD) measures the dissimilarity between word distributions and is defined as:

$$\text{JSD}(p||q) = \frac{1}{2} \left[ D_{\text{KL}}(p||m) + D_{\text{KL}}(q||m) \right] , \quad (4)$$

   where $D_{\text{KL}}$ is the Kulback-Leibler divergence, $m = (p+q)/2$. For TF-IDF compression, we calculate the JSD between the average word frequencies in the original and compressed documents, defined as:

$$p_i = \frac{1}{N} \sum_{j=1}^{N} f_{i,j} \, , \quad q_i = \frac{1}{N} \sum_{j=1}^{N} f'_{i,j} \, , \qquad (5)$$

   where $f_{i,j}$ and $f'_{i,j}$ are the frequencies of word $w_i$ in documents $d_j$ and $d'_j$, respectively.

   For LDA compression, we use the average JSD across all topics, defined as:

$$\overline{\text{JSD}}(p||q) = \sum_{t=1}^{T} \pi_t \, \text{JSD}(p_t, q_t) \, , \qquad (6)$$

$$\pi_t = \frac{1}{N} \sum_{j=1}^{N} z_{t,d_j} \, , \qquad (7)$$

   where $z_{t,d}$ is an indicator function that equals 1 if topic $t$ is the most probable topic for document $d$, and zero otherwise. The densities $p_t$ and $q_t$ are defined as:

$$p_{t,i} = \frac{1}{N_t} \sum_{j=1}^{N} f_{i,j} \, z_{t,d_j} \, , \qquad (8)$$

$$q_{t,i} = \frac{1}{N_t} \sum_{j=1}^{N} f'_{i,j} \, z_{t,d_j} \, , \qquad (9)$$

with $f_{i,j}$ and $f'_{i,j}$ given in (5), and $N_t$ is the number of documents for which topic $t$ is the most probable one. Further details on the properties of JSD are available in Lin (1991). This metric allows us to evaluate how well the compressed documents retain the original word distributions.

3. **ROUGE score.** As a summarization metric, used to evaluate the quality of text summarization, we use the ROUGE-LCS score, introduced in Lin (2004), where $\text{LCS}(r, s)$ denotes the length of the longest sequence of words that appear in both $r$ and $s$ in the same order. The ROUGE-F1 score is defined as:

$$\text{ROUGE-F1} = 2 \, \frac{R \, P}{R + P} \, , \qquad (10)$$

   where recall $R = |\text{LCS}(r, s)|/|r|$ and precision $P = |\text{LCS}(r, s)|/|s|$; $|r|$, $|s|$, and $|\text{LCS}(r, s)|$ are the word counts in the corresponding sequences. This metric is used to assess how well the compressed documents retain the key information of the original documents.

## 2.2 Encoding procedure

We now describe the steps of the encoding procedure, following the work by Kanerva (1988):

1. We consider the English alphabet plus digits, denoted as $\mathcal{A}$, and assign to each element $a_k \in \mathcal{A}$ a random vector $\phi(a_k)$ from the space $\mathcal{H} = \{\pm 1\}^D$, where $D$ is the dimension of the space. In this vector space, we define a coordinate-wise multiplication operation $\otimes$ and a coordinate-wise sign operation $\oplus$. The multiplication is a simple coordinate-wise product, while the sign operation is applied after a coordinate-wise summation, with the sign of zero defined as 1.

2. We use word-wise encoding. To encode a word, we apply a permutation operation $\rho$ to each character's vector $\phi(a_k)$, shifting all but the first coordinate to the left. The encoding vector for word $w_i$ is then:

$$\phi(w_i) = \bigotimes_{0 \le k < |w_i|} \rho^k(\phi(a_k)) \, , \qquad (11)$$

where $|w_i|$ is the number of characters in word $w_i$.

3. The document encoding is obtained by applying the sign operation to the coordinate-wise summation of all word vectors:

$$\phi(d) = \bigoplus_{i=1}^{|d|} \phi(w_i)\,. \tag{12}$$

The outcome of this encoding procedure is a function $\phi(d)$ that maps a text to the vector space $\mathcal{H}$.

# 3 Theoretical analysis

We divide our theoretical analysis into two main components: *compression* and *encoding*, based on the compression HDC model (Figure 1) and the previous section. These components are supported by intuition, assumptions and theorems in the following subsections.

## 3.1 Compression

In this section, we present our compression analysis separately for TF-IDF and LDA-based approaches. The original TF-IDF and LDA statistics were introduced by (Aizawa, 2003) and (Blei et al., 2003), respectively.

### 3.1.1 TF-IDF part

We analyze the TF-IDF score $\mathrm{ts}(w_i, \mathcal{D})$ as a random variable. The randomness stems from the frequency $f_{i,j}$ and the number of documents $N_{w_i}$ containing the word $w_i$. The frequency $f_{i,j}$ is related to the number of occurrences $n_{i,j}$ of word $w_i$ in document $d_j$ as $n_{i,j} = f_{i,j}|d_j|$. We can represent the documents schematically as:

$$d_j = \underbrace{w_1 \ldots w_1}_{n_{1,j}} \ldots \underbrace{w_M \ldots w_M}_{n_{M,j}}\,. \tag{13}$$

Thus, each document can be considered as a random vector $(n_{1,j}, n_{2,j}, ..., n_{M,j})$. To proceed with our analysis, we make the following assumptions:

**Assumption 1** (**Poisson-like distribution and independence across documents**). *To model the TF-IDF distribution, we assume that the number of occurrences $n_{i,j}$ of word $w_i$ in document $d_j$ are independent of the document $d_j$ and follows a distribution $Dist(\lambda_i)$, where:*

$$\mathbb{P}(n_{i,j} = k) = \begin{cases} 1 - f(\lambda_i)\,, & k = 0\,; \\ f(\lambda_i)\dfrac{\lambda_i^k e^{-\lambda_i}}{k!(1 - e^{-\lambda_i})}\,, & k > 0\,. \end{cases} \tag{14}$$

*Here, $f(\lambda_i)$ is an auxiliary function introduced to make our theoretical analysis tractable and ensure a monotonically growing TF-IDF approximate estimate, prioritizing words with larger $\lambda_i$ for encoding.*

The next assumption allows us to exclude randomness from the TF part:

**Assumption 2** (**Average frequency**). *The TF part can be fixed at $p_i$, by approximating the average frequency as:*

$$\frac{1}{N} \sum_{j=1}^{N} f_{i,j} = \frac{1}{N} \sum_{j=1}^{N} \frac{n_{i,j}}{|d_j|} \approx \frac{\mathbb{E}n_i}{\mathbb{E}|d_j|} = p_i\,, \tag{15}$$

*where $\mathbb{E}|d_j| = \displaystyle\sum_{i=1}^{M} \lambda_i.$*

Thus randomness retains only in the IDF part, i.e. in $N_w$. To estimate the number of documents where word $w$ occurs at least once, we have:

$$N_w = \sum_{j=1}^{N} \mathbf{1}(w \in d_j)\,, \tag{16}$$

which is a sum of $N$ i.i.d. Bernoulli variables $\mathrm{Bern}(q_w)$ with $q_w = 1 - \exp(-\lambda_w)$. Hence, the expectation of $N_w$ is $q_w N$, and for the TF-IDF approximate we obtain:

$$\widetilde{\mathrm{ts}}(w) = -\frac{\lambda_w f(\lambda_w)}{(1 - e^{-\lambda_w})\mathbb{E}|d|} \ln(1 - e^{-\lambda_w})\,. \tag{17}$$

To ensure a monotonically growing TF-IDF approximation, we make the next assumption:

**Assumption 3** (**Function $f(x)$**). *Function $f(x)$ is defined as:*

$$f(\lambda) = \frac{\lambda}{1 + \lambda}(1 - e^{-\lambda})\,. \tag{18}$$

This results in the following score approximate expectation:

$$\widetilde{\mathrm{ts}}(w) = -\frac{\lambda_w^2}{(1 + \lambda_w)\mathbb{E}|d|} \ln\left[\frac{\lambda_w}{1 + \lambda_w}(1 - e^{-\lambda_w})\right] \tag{19}$$

with the asymptotic behavior $\widetilde{\mathrm{ts}}(w)\mathbb{E}|d| = 1 - 3/(2\lambda) + \mathcal{O}(\lambda^{-2})$, i.e. attaining gradually 1 from below.

Figure 2 illustrates the true TF-IDF score (1) for IMDB dataset and our approximate expectation $\widetilde{\mathrm{ts}}(w)$ as a function of the parameter estimate $\widehat{\lambda}_w$, obtained using the method of moments from the equation:

$$n_w \equiv \frac{1}{N} \sum_{j=1}^{N} n_{w,j} = \frac{\widehat{\lambda}_w f(\widehat{\lambda}_w)}{1 - e^{-\widehat{\lambda}_w}} \tag{20}$$

(here and below, estimators of random variables are denoted with a wide hat). As can be observed, $\widetilde{\mathrm{ts}}(w)$

31

Figure 2: Comparison of the true TF-IDF statistics $\text{ts}(w)$ (1) for IMDB dataset and its approximate expectation $\widetilde{\text{ts}}(w)$ (19).

grows monotonically, as does the average true TF-IDF. However, the true TF-IDF values exhibit a noticeable vertical scatter (see blue points at $\lambda_w < 0.4$) due to the inherent randomness of the true TF-IDF score.

The compression method outlined in Section 2.1 selects words with the largest TF-IDF score:

$$\mathcal{W}_p = \{w \in \mathcal{W} : \text{ts}(w) \geq \text{ts}_{(\lceil (1-p)M \rceil)}\}. \quad (21)$$

Here and below $X_{(k)}$ denote is the $k$-th order statistic of $\{X(w_1), \ldots, X(w_M)\}$. Due to the complexity of $\text{ts}(w)$, we use expectation $\widetilde{\text{ts}}(w)$ to select the $pM$ words with the highest values of $\widehat{\lambda}_w$:

$$\widehat{\mathcal{W}}_p = \{w_i \in \mathcal{W} : \widehat{\lambda}_i \geq \widehat{\lambda}_*\}, \quad (22)$$

where $\widehat{\lambda}_* \equiv \widehat{\lambda}_{(\lceil (1-p)M \rceil)}$ is the minimal value $\widehat{\lambda}_w$ of the word $w$ included in set $\widehat{\mathcal{W}}_p$. Although $\mathcal{W}_p$ and $\widehat{\mathcal{W}}_p$ are not identical due to the randomness of $\text{ts}(w_i)$ and $\widehat{\lambda}_w$, the monotonicity of $\widetilde{\text{ts}}(w)$ implies that both sets will contain the same words, except for those in the vicinity of $\widehat{\lambda}_*$, where some words will be randomly added and others excluded from $\mathcal{W}_p$. To simplify our analysis, we assume that the sets $\mathcal{W}_p$ and $\widehat{\mathcal{W}}_p$ differ negligibly:

**Assumption 4** (**Negligible difference in selected words**). *We assume that $\mathcal{W}_p$ and $\widehat{\mathcal{W}}_p$ differ negligibly.*

For the theorems, we require an informational inequality (proof follows from Pinsker's inequality and Lin, 1991):

**Lemma 1.** *For Jensen-Shannon divergence, we have:*

$$\frac{1}{4}\left[V^2(p,m) + V^2(q,m)\right] \leq \text{JSD}(p||q) \leq \frac{1}{2}V(p,q), \quad (23)$$

*where $V(p,q) = \sum_i |p_i - q_i|$ and $m = (p+q)/2$.*

We now formulate the theorems (see Appendix A.1.1 for the proof).

**Theorem 1** (TF-IDF compression). *Based on assumptions 1–4, we have the consistent estimators for $\text{CR}, \text{JSD}(p||q)$ and ROUGE-F1:*

$$\widehat{\text{CR}} = \frac{\sum_{w \in \widehat{\mathcal{W}}_p} g(\widehat{\lambda}_w)}{\sum_{w \in \mathcal{W}} g(\widehat{\lambda}_w)}, \quad (24)$$

$$\widehat{\text{JSD}}(p||q) = \frac{1}{2}\left[\sum_{w \in \widehat{\mathcal{W}}_p} \widehat{p}_w \ln\left(\frac{2\widehat{\text{CR}}}{\widehat{\text{CR}} + 1}\right)\right]$$

$$+ \frac{\ln 2}{2}\sum_{w \in \mathcal{W}/\widehat{\mathcal{W}}_p} \widehat{p}_w + \frac{1}{2}\left[\sum_{w \in \widehat{\mathcal{W}}_p} \frac{\widehat{p}_w}{\widehat{\text{CR}}} \ln\left(\frac{2}{1 + \widehat{\text{CR}}}\right)\right], \quad (25)$$

$$\widehat{\text{ROUGE-F1}} = 2\frac{\widehat{\text{CR}}}{\widehat{\text{CR}} + 1}, \quad (26)$$

*where $g(x) = x^2/(1 + x)$ and $\widehat{p}_w = g(\widehat{\lambda}_w)/\sum_{w \in \mathcal{W}} g(\widehat{\lambda}_w)$.*

**Theorem 2** (Quantile criteria). *Under assumptions 1–4, the TF-IDF compression model with $p$-quantile criteria has the following bounds from Table 1.*

### 3.1.2 LDA part

We now examine the LDA compression procedure. For a fixed topic $t$, the distribution of words is a Dirichlet random variable, $\Phi_t \sim Dir(\alpha)$, where $\alpha$ is a vector of parameters $(\alpha_1, \ldots, \alpha_M)$ (see Blei et al., 2003, for details). As outlined in Section 2.1, we define the set:

$$\mathcal{W}_{t,p} = \{w_i \in \mathcal{W} : \Phi_{t,w} \geq \Phi_{t,(\lceil (1-p)M \rceil)}\}, \quad (27)$$

where $\Phi_{t,w}$ is the probability of word $w$ belonging to topic $t$. To determine the distribution of $\Phi_{t,(\lceil (1-p)M \rceil)}$, we need the marginal distributions of $\Phi_{t,w_i}$.

**Lemma 2.** *If $\Phi = (\Phi_1, \ldots, \Phi_M) \sim Dir(\alpha)$, then its marginal distributions are beta distributed random variables:*

$$\Phi_i \sim Beta\left(\alpha_i, \sum_{k=1}^{M} \alpha_k - \alpha_i\right). \quad (28)$$

This lemma allows us to identify and generalize the object of our interest. Applying the same conceptual approach as in the TF-IDF part, we focus on the quantile value of the $(\Phi_{t,1}, \ldots, \Phi_{t,M})$, where each $\Phi_{t,i}$ is distributed as in (28).

The model has an additional parameter $\alpha$, which we set to $(0.5, \ldots, 0.5)$, implying that we are unsure about word significance in topic $t$:

32

**Assumption 5 (Non-significance).** $\alpha = (0.5, \ldots, 0.5)$.

Under Assumption 5, we have a set of $Beta(0.5, 0.5[M-1])$ random variables. Using the same expectation approach as in the TF-IDF case, we focus on estimating $\mathbb{E}\Phi_{t,(k)}$. To proceed, we use the following lemma (see Arnold and Groeneveld, 1979, for the proof):

**Lemma 3.** *For i.i.d. random variables $X_1, \ldots, X_n$ with mean $\mu$ and variance $\sigma^2$, we have the following inequality:*

$$-\sigma\sqrt{\frac{n-k}{k}} \leq \mathbb{E}X_{(k)} - \mu \leq \sigma\sqrt{\frac{k-1}{n-k+1}}. \quad (29)$$

For $X \sim Beta(\alpha, \beta)$, we have:

$$\mu = \frac{\alpha}{\alpha+\beta} = M^{-1}, \quad (30)$$

$$\sigma^2 = \frac{\alpha\beta}{(\alpha+\beta)^2(\alpha+\beta+1)} = \quad (31)$$

$$\frac{M-1}{M^2(0.5M+1)} \approx 2M^{-2}. \quad (32)$$

Hence, we can estimate the bounds of $\mathbb{E}\Phi_{t,(\lceil(1-p)M\rceil)}$.

Before proceeding with the theorems, we clarify the distribution of the number of occurrences. Unlike the TF-IDF model, where we calculated $n_{i,j}$ directly, in the LDA model, we operate with $\Phi_{t,i}$ values. Therefore, we assume:

**Assumption 6 (Poisson distribution).** *For each topic $t$, we assume that the number of occurrences of each word $w_i$ in a document $d_t$ are independent random variables following the Poisson distribution:*

$$d_t = \underbrace{w_1 \ldots w_1}_{v_{t,1} \sim Pois(\Phi_{t,1}C)} \quad \ldots \quad \underbrace{w_M \ldots w_M}_{v_{t,M} \sim Pois(\Phi_{t,M}C)}, \quad (33)$$

*where $v_{t,i}$ is the number of occurrences of word $w_i$ in a document $d_t$ belonging to topic $t$.*

This assumption is quite strict, as it assumes a constant $C$ that regulates the number of occurrences of each word in the document, and that this constant is the same for all topics. As we argue below, we use it to estimate the number of words in a document on a given topic.

Given a matrix of words in topic probabilities $\widehat{\Phi}_{t,w}$, we formulate the following theorems:

**Theorem 3 (LDA compression estimators).** *Under assumptions 5–6, we have asymptotically-unbiased estimators for* CR, JSD$(p||q)$, *and*

ROUGE-F1*:*

$$\widehat{\text{CR}} = \frac{\sum_{t=1}^{T} \pi_t \sum_{w \in \mathcal{W}_{p,t}} \widehat{\Phi}_{t,w}}{\sum_{t=1}^{T} \pi_t \sum_{w \in \mathcal{W}} \widehat{\Phi}_{t,w}}, \quad (34)$$

$$\widehat{\text{JSD}}(p||q) = \frac{1}{2}\sum_{t=1}^{T} \pi_t \left[\sum_{w \in \mathcal{W}} \widehat{\Phi}_{t,w} \ln\left(\frac{2\widehat{\text{CR}}}{\widehat{\text{CR}}+1}\right)\right]$$

$$+ \frac{1}{2}\sum_{t=1}^{T} \pi_t \left[\sum_{w \in \mathcal{W}} \frac{\widehat{\Phi}_{t,w}}{\widehat{\text{CR}}} \ln\left(\frac{2}{1+\widehat{\text{CR}}}\right)\right]$$

$$+ \frac{\ln 2}{2}\sum_{t=1}^{T} \pi_t \sum_{w \in \mathcal{W}\backslash\mathcal{W}_p} \widehat{\Phi}_{t,w}, \quad (35)$$

$$\widehat{\text{ROUGE}}\text{-F1} = 2\frac{\widehat{\text{CR}}}{\widehat{\text{CR}}+1} \quad (36)$$

*with $\pi_t$ defined by Eq. (7).*

**Theorem 4 (LDA compression bounds).** *Under assumptions 5–6, the LDA compression model with p-quantile criteria has the following bounds from Table 1.*

### 3.2 Encoding

To prove the applicability of our proposed CHDC approach, we now turn to encoding implications and focus on estimating the quality of document analysis based on an average document size. As in the previous section, we consider documents as a bag of words (13). Consider now two documents, $d_1$ and $d_2$. Given the binary HDC encoding, we map our documents to the $\phi(d_1)$ and $\phi(d_2)$, according to the rules from Section 2.2. As pointed in (Kanerva, 1988), the HDC model should distinguish the vectors $\phi(d_1)$ and $\phi(d_2)$, which means that:

$$\langle\phi(d_1), \phi(d_2)\rangle \to 0 \quad (37)$$

with $D \to \infty$ (here $\langle., .\rangle$ denotes the standard Euclidean dot-product). To estimate the effect of the encoding under fixed $D$, we propose to consider:

$$\mathbb{P}(\langle\phi(d_1), \phi(d_2)\rangle \geq \varepsilon D), \quad (38)$$

where $D$ is the vector-space dimension, $\varepsilon$ is small parameter that characterize distinguishability, $\phi$ is the encoding function, mentioned before. Notice that the $\phi(d)$ is a random vector, since we use a random binary HDC encoding. Therefore, we need to be sure that the probability of $\mathbb{P}(\langle\phi(d_1), \phi(d_2)\rangle > \varepsilon D)$ is low.

Let's rewrite the dot-product as follows:

$$\langle\phi(d_1), \phi(d_2)\rangle = \sum_{i=1}^{D} \phi_{1,i}\phi_{2,i} = \sum_{i=1}^{D} X_i, \quad (39)$$

33

| Th. | CR | JSD | ROUGE-F1 |
|---|---|---|---|
| Th. 2 | $\left[pM\dfrac{\min_{w\in\mathcal{W}_p}g(\hat{\lambda}_w)}{\sum_{w\in\mathcal{W}}g(\hat{\lambda}_w)};pM\dfrac{\max_{w\in\mathcal{W}_p}g(\hat{\lambda}_w)}{\sum_{w\in\mathcal{W}}g(\hat{\lambda}_w)}\right]$ | $\left[\dfrac{1}{4}\left[\widehat{V}_{pm}^2+\widehat{V}_{qm}^2\right];\dfrac{1}{2}\widehat{V}_{pq}\right]$ | $\left[2\dfrac{\mathrm{CR}_{\min}}{1+\mathrm{CR}_{\min}};2\dfrac{\mathrm{CR}_{\max}}{1+\mathrm{CR}_{\max}}\right]$ |
| Th. 4 | $\left[p-\sqrt{\dfrac{2p}{1-p}};p+p\sqrt{2(M-1)}\right]$ | $\left[\dfrac{1}{4}\sum_t\widehat{\pi}_t\left[\widehat{V}_{t,pm}^2+\widehat{V}_{t,qm}^2\right];\dfrac{1}{2}\sum_t\widehat{\pi}_t\widehat{V}_{t,pq}\right]$ | $\left[2\dfrac{\mathrm{CR}_{\min}}{1+\mathrm{CR}_{\min}};2\dfrac{\mathrm{CR}_{\max}}{1+\mathrm{CR}_{\max}}\right]$ |

Table 1: Bounds for the performance metrics: compression rate (CR), Jensen-Shannon divergence (JSD), and ROUGE-F1 score, under TF-IDF (Theorem 2) and LDA (Theorem 4) compression.

where $X_i$ are dependent Bernoulli-type random variables taking values in $\{\pm 1\}$, with $\gamma_i(X_1,\ldots,X_{i-1},\,X_{i+1},\ldots,X_D) = \mathbb{P}(X_i = 1|\{X_1,\ldots,X_D\}\setminus X_i)$. Unfortunately, we can't directly apply known techniques due to the possible dependency of the $\{X_i\}_{i=1}^D$. However, we propose the following lemma to overcome this problem (for proof, see Appendix A.2):

**Lemma 4.** *Assume $\{X_i\}_i$ are dependent random variables with Bernoulli-type distribution and $\mathbb{P}(X_i = 1|X_{i_1},\ldots,X_{i_k}) \leq p$. Then there are $\{Y_i\}_i$ independent Bernoulli variables with $\mathbb{P}(Y_i = 1) = p$ and we have:*

$$\mathbb{P}\left(\sum_{i=1}^D X_i \geq \varepsilon D\right) \leq \mathbb{P}\left(\sum_{i=1}^D Y_i \geq \varepsilon D\right). \quad (40)$$

The given lemma allows us to consider $X_i$ as independent random variables with the same bound $\gamma$ on its probability. To estimate the value of probability in (38), we propose using the following lemma (see Chernoff, 1952, for proof):

**Lemma 5** (Chernoff bound). *For a sum of independent random variables $X = \sum_i X_i$, we have:*

$$\mathbb{P}(X \geq a) \leq \inf_{t>0}\left[e^{-ta}\prod_i \mathbb{E}e^{tX_i}\right]. \quad (41)$$

To justify the model, we formulate the following theorem (for the proof, see Appendix A.2):

**Theorem 5.** *The probability (38) is upper bounded by:*

$$\mathbb{P}(\langle\phi(d_1),\phi(d_2)\rangle \geq \varepsilon D) \leq F(D,\gamma,\varepsilon), \quad (42)$$

*where:*

*1. The upper boundary:*

$$\ln F(D,\varepsilon,\gamma) = \frac{D}{2}(1-\varepsilon)\ln\left[\frac{1-\gamma}{\gamma}\frac{1+\varepsilon}{1-\varepsilon}\right]$$
$$- D\ln\left[\frac{1+\varepsilon}{2\gamma}\right]. \quad (43)$$

*2. The Bernoulli probability $\gamma$ satisfies the inequality:*

$$\frac{1}{2} < \gamma \leq \frac{1}{2} + \binom{|d|}{\lceil|d|/2\rceil}\frac{1}{2^{|d|}} \approx \frac{1}{2} + \sqrt{\frac{2}{\pi|d|}}, \quad (44)$$

*where $|d| = \mathbb{E}|d_i|$ is the average length of the document, the round brackets denote the binomial coefficient, and the asymptotical expansion in the r.h.s is obtained using Stirling's approximation.*

The function $F$ attains a maximum value of 1 when $\varepsilon + 1 = 2\gamma$. As we move away from this line, the function rapidly declines, with the decline becoming sharper as $D$ increases. This implies that

$$\varepsilon \lesssim 2\sqrt{\frac{2}{\pi|d|}}. \quad (45)$$

For example, in the IMDB dataset, compression for $p = 0.1$ from an average document length of 122 words to 100 words increases $\varepsilon$ by a factor of approximately $\sqrt{122/100} \approx 1.1$, just slightly worsening distinguishability.

## 4 Experiments

To verify our theoretical results, we propose a two-stage experimental setup, focusing on compression effect estimation and encoding results.

### 4.1 Compression analysis

We explore TF-IDF and LDA text compression techniques using Algorithm 1 (see A.4) applying it to IMDB reviews (Maas et al., 2011), AG News Dataset (Zhang et al., 2015), and arXiv dataset (Clement et al., 2019). Figure 3 (see A.3) presents the results, comparing direct calculations of the three metrics (CR, JSD and ROUGE-F1) with their theoretical expectations for different quantile parameters $p$. The green bounds show the possible ranges of metric scatter due to the randomness of word distributions (Theorems 2 and 4). The three upper panel rows demonstrate that TF-IDF

| | TF-IDF | | | LDA | | |
|---|---|---|---|---|---|---|
| $D$ | $\widehat{\varepsilon}_{p=0.01}$ | $\widehat{\varepsilon}_{p=0.1}$ | $\widehat{\varepsilon}_{p=1}$ | $\widehat{\varepsilon}_{p=0.01}$ | $\widehat{\varepsilon}_{p=0.1}$ | $\widehat{\varepsilon}_{p=1}$ |
| 256 | $0.17 \pm 0.02$ | $0.13 \pm 0.01$ | $0.12 \pm 0.01$ | $0.16 \pm 0.02$ | $0.13 \pm 0.01$ | $0.12 \pm 0.01$ |
| 1024 | $0.17 \pm 0.02$ | $0.13 \pm 0.01$ | $0.12 \pm 0.01$ | $0.16 \pm 0.01$ | $0.12 \pm 0.01$ | $0.12 \pm 0.01$ |
| 4096 | $0.17 \pm 0.02$ | $0.13 \pm 0.01$ | $0.12 \pm 0.01$ | $0.16 \pm 0.01$ | $0.12 \pm 0.01$ | $0.12 \pm 0.01$ |
| 16384 | $0.17 \pm 0.02$ | $0.12 \pm 0.01$ | $0.11 \pm 0.01$ | $0.16 \pm 0.01$ | $0.12 \pm 0.01$ | $0.11 \pm 0.01$ |

Table 2: Encoding analysis for TF-IDF and LDA compression techniques using the IMDB dataset. The table shows average scalar product values for dictionary compression parameters $p = 0.01$, $0.1$, and $1$ ($|d| \approx 60$, $100$, $122$, respectively) and vector space dimension $D$.

compression accurately captures all three metrics across all datasets and different values of $p$, because the relevant variables are directly observed and the assumptions are reasonable. In contrast, the three lower panels show that the LDA compression estimators perform worse, likely because the underlying distributional assumptions do not fully correspond to the actual distributions.

## 4.2 Encoding analysis

To validate the results in Section 3.2, we analyze how the encoding procedure impacts the distinguishability of randomly selected documents using the IMDB dataset. This dataset, which comprises two classes, simplifies our analysis (Algorithm 2) while still revealing key insights. We use Monte Carlo simulations with 100 iterations for the alphabet $\mathcal{A}$ and 100 iterations for document sampling (pairs from different classes), resulting in 10000 total iterations. Table 2 presents estimates of the parameter $\varepsilon$, defined as:

$$\widehat{\varepsilon}_p = D^{-1} \, \mathbb{E} |\langle \phi(d_{1,p}), \phi(d_{2,p}) \rangle| \qquad (46)$$

where $d_{1,p}$ and $d_{2,p}$ are randomly selected documents from different classes after compression, and $p$ is the compression parameter. The table shows results for $p = 1$ (no compression, $|d| \approx 122$ words), $p = 0.1$ (medium compression, $|d| \approx 100$ words), and $p = 0.01$ (high compression, $|d| \approx 60$ words).

The estimates $\widehat{\varepsilon}_p$ are similar for TF-IDF and LDA compression techniques, decreasing approximately with the square root of the average document size $|d|$ and remaining within 20% of the theoretical upper boundary (45).

## 5 Discussion

This paper introduces a novel approach to address dimensionality concerns in Hyperdimensional Computing (HDC) by adding compression. We propose a model that combines TF-IDF or LDA-based compression with binary HDC to mitigate the curse of dimensionality. Section 3 presents the

core concepts, and Section 4 provides experimental results validating our approach. Our method demonstrates that significantly reducing the encoding space of the initial dictionary only slightly compromises class distinguishability in classification tasks. Specifically, reducing the dictionary by 10 times increases the distinguishability parameter by 10%, and reducing it by 100 times increases the parameter by 40%, while still maintaining a low value (far from 1).

Theorems 1 and 3 provide estimators that accurately estimate the necessary parameters, with TF-IDF compression showing particularly low error and LDA offering slightly better explainability in encoding analysis (see Table 2).

Despite our numerical results aligning with theory, we identify two drawbacks that warrant further research and development in this field:

1. We observe that the bounds provided in Theorems 2–4 are not sufficiently tight. Because these bounds are estimated using the distribution properties of the datasets, it is difficult to obtain tighter bounds for the given metrics.

2. The encoding effect diminishes with increasing vector space size $D$. This effect, explained by Theorem 3.2, is due to the upper boundary function $F$ becoming concentrated in a narrow region near the line $\varepsilon + 1 = 2\gamma$ as $D$ increases, which reduces the confidence intervals of the estimates $\widehat{\varepsilon}$, without lowering the estimates themselves.

Our results provide several insights into the application of TF-IDF- and LDA-based compression techniques and demonstrate the potential of Compression HDC for broader practical application to empirical problems, where noise significantly hinders data compression and classification.

# References

Akiko Aizawa. 2003. An information-theoretic perspective of tf—idf measures. *Inf. Process. Manage.*, 39(1):45–65.

Barry C. Arnold and Richard A. Groeneveld. 1979. Bounds on expectations of linear systematic statistics based on dependent samples. *Annals of Statistics*, 7:220–223.

David M. Blei, Andrew Y. Ng, and Michael I. Jordan. 2003. Latent dirichlet allocation. *Journal of Machine Learning Research*, 3:993–1022.

Herman Chernoff. 1952. A Measure of Asymptotic Efficiency for Tests of a Hypothesis Based on the sum of Observations. *The Annals of Mathematical Statistics*, 23(4):493 – 507.

Colin B. Clement, Matthew Bierbaum, Kevin P. O'Keeffe, and Alexander A. Alemi. 2019. On the use of arxiv as a dataset.

P Kanerva. 1988. *Sparse Distributed Memory*. MIT Press, Cambridge, MA.

Pentti Kanerva. 2009. Hyperdimensional computing: An introduction to computing in distributed representation with high-dimensional random vectors. *Cognitive Computation*, 1(2):139–159.

Denis Kleyko, Dmitri Rachkovskij, Evgeny Osipov, and Abbas Rahimi. 2023. A survey on hyperdimensional computing aka vector symbolic architectures, part ii: Applications, cognitive models, and challenges. *ACM Comput. Surv.*, 55(9).

Chin-Yew Lin. 2004. Rouge: A package for automatic evaluation of summaries. In *Proceedings of the Workshop on Text Summarization Branches Out (WAS 2004)*, pages 74–81. Association for Computational Linguistics.

Jianhua Lin. 1991. Divergence measures based on the shannon entropy. *IEEE Transactions on Information Theory*, 37(1):145–151.

Andrew L. Maas, Raymond E. Daly, Peter T. Pham, Dan Huang, Andrew Y. Ng, and Christopher Potts. 2011. Learning word vectors for sentiment analysis. In *Proceedings of the 49th Annual Meeting of the Association for Computational Linguistics: Human Language Technologies*, pages 142–150, Portland, Oregon, USA. Association for Computational Linguistics.

Anton Mitrokhin, Peter Sutor, Cornelia Fermüller, and Yiannis Aloimonos. 2019. Learning sensorimotor control with neuromorphic sensors: Toward hyperdimensional active perception. *Science Robotics*, 4:eaaw6736.

Peer Neubert and Stefan Schubert. 2021. Hyperdimensional computing as a framework for systematic aggregation of image descriptors. pages 16933–16942.

Saeid Pourmand, Wyatt D. Whiting, Alireza Aghasi, and Nicholas F. Marshall. 2024. Laplace-hdc: Understanding the geometry of binary hyperdimensional computing. *ArXiv*, abs/2404.10759.

Abbas Rahimi, Pentti Kanerva, and Jan M. Rabaey. 2019. Efficient biosignal processing using hyperdimensional computing: A case study for emg-based hand gesture recognition. *IEEE Transactions on Biomedical Engineering*, 66(11):3192–3203.

Karen Spärck Jones. 1972. A statistical interpretation of term specificity and its application in retrieval. *Journal of Documentation*, 28(1):11–21.

Anthony Thomas, Sanjoy Dasgupta, and Tajana Rosing. 2021. A theoretical perspective on hyperdimensional computing. *Journal of Artificial Intelligence Research*, 72:215–249.

Tao Yu, Yichi Zhang, Zhiru Zhang, and Christopher De Sa. 2024. Understanding hyperdimensional computing for parallel single-pass learning. In *Proceedings of the 36th International Conference on Neural Information Processing Systems*, NIPS '22, Red Hook, NY, USA. Curran Associates Inc.

Guo Yunhui, Mohsen Imani, Jaeyoung Kang, Sahand Salamat, Justin Morris, Baris Aksanli, Yeseong Kim, and Tajana Rosing. 2021. Hyperrec: Efficient recommender systems with hyperdimensional computing. pages 384–389.

Xiang Zhang, Junbo Jake Zhao, and Yann LeCun. 2015. Character-level convolutional networks for text classification. In *NIPS*.

## A  Appendix / supplemental material

### A.1  Compression analysis

In the given section, we provide the theoretical justification of the analysis provided in the paper before. The first part of the upcoming appendix correspond to the TF-IDF and LDA theories.

#### A.1.1  TF-IDF part

*Lemma 1.* From Theorem 3 in (Lin, 1991) we have:

$$\text{JSD}(p||q) \leq \frac{1}{2}V(p,q)$$

Using definition of JSD and Pinsker inequality:

$$\text{JSD}(p||q) = \frac{1}{2}\left[D_{KL}(p||m) + D_{KL}(q||m)\right] \geq$$

$$\frac{1}{4}\left[V(p,m) + V(q,m)\right]$$

Now, we are ready to move to the proofs of the theorems.

*Theorem 1-Theorem 2.* 1. Follow the definition of CR, we have:

$$\text{CR} = \frac{\frac{1}{N}\sum_j |d'|_j}{\frac{1}{N}\sum_j |d|_j} \simeq \frac{\mathbb{E}|d'|}{\mathbb{E}|d|}$$

Based on the model in Assumption 1 we have:

$$\mathbb{E}|d| = \sum_{w \in \mathcal{W}} \frac{\lambda_w^2}{1 + \lambda_w} \tag{47}$$

Notice that after the compression procedure, we leave only the words from $\widehat{\mathcal{W}}_p$; hence, given the (47), we have:

$$\text{CR} \simeq \frac{\sum_{w \in \mathcal{W}_p} g(\lambda_w)}{\sum_{w \in \mathcal{W}} g(\lambda_w)},$$

where $g(x) = \dfrac{x^2}{1 + x}$. We obtain the result of the Theorem 1.1 by using the consistent estimator (20) for $\lambda_w$ and using Slutsky's theorem.

Also, we easily obtain the bounds for Theorem 2.1 for $\widehat{\text{CR}}$:

$$\left(pM\frac{\min_{w \in \mathcal{W}_p} g(\widehat{\lambda}_w)}{\sum_w g(\widehat{\lambda}_w)}, pM\frac{\max_{w \in \mathcal{W}_p} g(\widehat{\lambda}_w)}{\sum_w g(\widehat{\lambda}_w)}\right) \tag{48}$$

2. Using the Jensen-Shannon divergence definition and Lemma 3 we have:

$$\text{JSD}(p||q) = \frac{1}{2}\left[D_{KL}(p||m) + D_{KL}(q||m)\right],$$

where $p = \{p_w\}$ and $q = \{q_w\}$, defined in (5). Notice that based on Assumption 1we have the following form for $p_w$ and $q_w$:

$$p_w = \frac{n_w}{|d|}, \quad q_w = \frac{n_w}{|d'|} \tag{49}$$

Hence we have $\text{CR}p_w = q_w$. Next, we can easily find the consistent estimator for $p_w$:

$$\widehat{p}_w = \frac{\widehat{\lambda}_w}{\sum_k g(\widehat{\lambda}_k)}, \tag{50}$$

because of Slutsky's theorem and consistent estimator for $\lambda_w$. Now, using the definition of $D_{KL}$:

$$D_{KL}(p||q) = \sum_w p_w \log\left(\frac{p_w}{q_w}\right), \tag{51}$$

and previous properties: $\text{CR} \times p_w = q_w$ and $q_w = 0$ for $w \in \mathcal{W} \setminus \widehat{\mathcal{W}}_p$ we obtain the results.

For the bounds in Theorem 2 we use Lemma 3.

3. ROUGE-L score. Here, we focus on the classical text compression score. ROUGE-L has three components to analyze:

1. Precision: $P = \dfrac{\mathbb{E}|LCS|}{\mathbb{E}_q|d|}$

2. Recall: $R = \dfrac{\mathbb{E}|LCS|}{\mathbb{E}_p|d|}$

3. F-score: $F1 = 2\dfrac{R \cdot P}{R + P}$

Notice that our procedure preserves the order, hence $\mathbb{E}|LCS| = \mathbb{E}_q|d|$. Hence, we have the following:

1. Precision: $P \equiv 1$

2. Recall: $R = \text{CR}$

3. F-score: $F1 = 2\dfrac{\text{CR}}{\text{CR} + 1}$

Now, since $f(x) = \dfrac{x}{x + 1}$ is increasing for $x \geq 0$, we proved our bounds.

#### A.1.2  LDA part

*Theorem 3-4.*

1. Notice that $\text{CR} \simeq \dfrac{\mathbb{E}|d'|}{\mathbb{E}|d|}$, hence using Assumption 6

$$\mathbb{E}|d| = \sum_{i=1}^{M} \mathbb{E}v_i = \sum_{i=1}^{M}\sum_{t=1}^{T} \pi_t C\Phi_{t,i},$$

where $\pi_t$ - probability of document's topic is $t$. Hence using the

$$\widehat{\text{CR}} = \frac{\sum_{i=1}^{M} \sum_{w \in \mathcal{W}_{t,tp}} \widehat{\pi}_t \Phi_{t,w}}{\sum_{i=1}^{M} \sum_w \widehat{\pi}_t \Phi_{t,w}},$$

where $\widehat{\pi}_t = \frac{1}{N} \sum_{j=1}^{N} z_{t,d_j}$ we obtain the consistent estimator of the CR.

The upper bound can be obtained as follows:

$$\widehat{\text{CR}} = \sum_t \pi_t \sum_{w \in \mathcal{W}_{p,t}} \Phi_{w,t},$$

where $\Phi_{w,t} \approx \mathbb{E} X_{(j)}$, $j$ corresponding number of order statistics and $X = \{X_1, \dots, X_M\}$ sequence of Beta distributed r.v. as in 2. Hence using the $\sum_t \pi_t = 1$, we can proceed with the Lemma 3 to obtain:

$$\widehat{\text{CR}} \geq p - \frac{\sqrt{2}}{M} \sum_{i=\lceil (1-p)M \rceil}^{M} \sqrt{\frac{M-i}{i}} \qquad (52)$$

$$\widehat{\text{CR}} \leq \left( p + \frac{\sqrt{2}}{M} \sum_{i=\lceil (1-p)M \rceil}^{M} \sqrt{\frac{i-1}{M-i+1}} \right) \qquad (53)$$

This leads us to the following:

$$p - \sqrt{\frac{2p}{1-p}} \leq \widehat{\text{CR}} \leq p + p\sqrt{2(M-1)}$$

2. We want to examine the value of the:

$$\overline{\text{JSD}(p||q)} = \sum_{t=1}^{T} \pi_t \text{JSD}(p_t, q_t),$$

where $p_{t,i} = \frac{1}{N_t} \sum_{j=1}^{N} f_{i,j} z_{t,d_j}$ and $q_{t,i} = \frac{1}{N_t} \sum_{j=1}^{N} f'_{i,j} z_{t,d_j}$. Under assumption Assumption 6, we have:

$$p_{t,i}/q_{t,i} = f_i/f'_i = 1/\text{CR}$$

Therefore, we have: $\text{CR} \times p_{t,i} = q_{t,i}$. Also, we have:

$$\widehat{p}_{t,i} = \frac{C \times \Phi_{t,i}}{\sum_k C \times \Phi_{t,k}} = \Phi_{t,i} \xrightarrow{\mathbb{P}} p_{t,i},$$

hence using Slutsky's theorem and consistent estimators for $\pi_t$ and $p_{t,i}, q_{t,i}$ we have the consistent estimator.

Bounds for JSD are obtained as in the proof of Theorem 2, using the definition (6)

3. The same idea as in the proof of the Theorem 2 works here.

## A.2 Encoding analysis

In the given section, we provide the theoretical justification of the encoding analysis, provided in the paper.

*Lemma 4.* Let's consider $u_1, \dots, u_D$ independent uniform distributions on $[0, 1]$. Denote $Y_i = \mathbf{1}(u_i \leq p)$, then $\{Y_i\}_i$ are independent. Here we assume $\mathbf{1}(...) \in \{\pm 1\}$, to satisfy the Bernoulli-type distribution of $X_i$.

Notice that $\mathbb{P}(X_i = 1) = \mathbb{P}(u_i \leq q_i)$, where $q_i = \mathbb{P}(X_i = 1 | X_1, \dots, X_{i-1})$ and thence:

$$X_i \leq Y_i \Rightarrow \mathbb{P}(\sum_{i=1}^{D} X_i \geq \varepsilon D) \leq \mathbb{P}\left( \sum_{i=1}^{D} Y_i \geq \varepsilon D \right)$$

*Theorem 5.*
**Probability estimation part.**
In the given appendix, we justify the ideas provided in the **encoding** part in the theory section. Notice that we aimed to consider the given probability:

$$\mathbb{P}(\langle \phi(d_1), \phi(d_2) \rangle \geq \varepsilon D) =$$
$$\mathbb{P}\left( \sum_{i=1}^{D} X_i \geq \varepsilon D \right) = \bigstar$$

Using the Lemma 5, we can obtain:

$$\bigstar \leq \inf_{t>0} \left[ e^{-\varepsilon Dt} \left( \mathbb{E} e^{tX} \right)^D \right],$$

where $X$ is a Bernoulli random variable with parameter $\gamma$ and values in $\{\pm 1\}$. Hence, we have:

$$\bigstar \leq \inf_{t>0} \left[ e^{-\varepsilon Dt} \left( \gamma e^t + (1-\gamma) e^{-t} \right)^D \right] = \inf_{t>0} L(t)$$

To find the minimum of the $L(t)$, we need to derive the first-order condition:

$$\frac{\mathrm{d}}{\mathrm{d}t} L(t) = 0$$

This is equivalent to:

$$\underbrace{(\gamma(e^{2t} - 1) + 1)^{D-1}}_{>0, \text{ since } \gamma < 1}$$
$$\times \left( (\gamma - 1)(\varepsilon D + D) - \gamma e^{2t}(\varepsilon D - D) \right) = 0$$

$$(1-\gamma)D(\varepsilon+1) = \gamma D(1-\varepsilon)e^{2t} \Rightarrow$$

$$t_{\min} = \frac{1}{2}\ln\left[\frac{1-\gamma}{\gamma}\frac{1+\varepsilon}{1-\varepsilon}\right] = \frac{1}{2}\ln\underbrace{C(\varepsilon)C(\gamma)}_{C(\varepsilon,\gamma)}$$

After rearranging, we have:

$$\exp\Big[-D\big(\varepsilon\ln\sqrt{C(\varepsilon,\gamma)}$$

$$-\ln\big(p\sqrt{C(\varepsilon,\gamma)}+\tfrac{1-\gamma}{\sqrt{C(\varepsilon,\gamma)}}\big)\big)\Big] =$$

$$\exp\left[-D\ln\left(\frac{C(\varepsilon,\gamma)^{(\varepsilon+1)/2}}{1-\gamma+\gamma C(\varepsilon,\gamma)}\right)\right] =$$

$$\exp\left[-D\ln\left(\underbrace{\frac{1}{2}\left[\frac{1-\gamma}{\gamma}\frac{1+\varepsilon}{1-\varepsilon}\right]^{(\varepsilon-1)/2}\frac{1+\varepsilon}{\gamma}}_{\star\star}\right)\right]$$

Hence, this probability decreases with increasing $D$ or by managing the expression in scopes. Simple algebra shows that for the same level of $D$ and $\varepsilon$, we can increase the expression $\star\star$ by increasing the $\gamma$ value after the critical point $\gamma_\varepsilon = \dfrac{1+\varepsilon}{2}$.

**Compression connection part.**

Next, we aim to connect the encoding analysis with the compression part. We provide the following explanation. Consider the following relationship:

$$\gamma = \mathbb{P}(\phi_{1,i}\phi_{2,i}=1) = \tilde{\gamma}^2 + (1-\tilde{\gamma})^2$$

where $\phi_i$ is the $i$-th position of the vector-encoding of randomly generated document $d$.

Notice that:

$$\tilde{\gamma} = \mathbb{P}\left(\text{sign}\left[\sum_{j=1}^{|d|}\phi_{i,j}\right]=1\right) =$$

$$\mathbb{P}\left(\text{sign}\left[\underbrace{\sum_{k=1}^{M}\#\{w_k\}\phi_{i,w_k}}_{\nu_i}\right]=1\right),$$

where the support of the $\nu_i$ is determined by the all possible sums of $\sum_{k=1}^{M}\pm\#\{w_k\}$. The behavior of

this sum is quite unpredictable, but we can say that the given distribution is symmetrical. To estimate $\mathbb{P}(\text{sign }\nu_i=1)$ we will consider the probability of $\eta = \mathbb{P}(\nu_i=0)$. Hence (by symmetry), we have:

$$\tilde{\gamma} = \frac{1}{2} + \frac{\eta}{2},$$

i.e., we cut half of the probability from the left tail of the distribution and add it to the right one. We propose the following estimation of the $\eta$:

$$\eta \leq \binom{|d|}{\lceil|d|/2\rceil}\frac{1}{2^{|d|}}$$

This bound is easy to obtain assuming $\nu_i \approx \sum_{i=1}^{|d|}\upsilon_i$, where $\upsilon_i$ is independent Bernoulli r.v. with values $\pm1$ and equal proabilities.

Based on the CR definition, $\text{CR}\times|d| = |d'|$, hence for compressed object the value of $\eta$ will be bounded by:

$$\eta \leq \binom{\text{CR}\,|d|}{\lceil\text{CR}\,|d|/2\rceil}\frac{1}{2^{\text{CR}\,|d|}}$$

The RHS is increasing with the decreasing of the CR. As a result, we have:

$$\gamma = \tilde{\gamma}^2 + (1-\tilde{\gamma})^2 = \left(\frac{1}{2}+\frac{\eta}{2}\right)^2 + \left(\frac{1}{2}-\frac{\eta}{2}\right)^2 \leq$$

$$\frac{1}{2} + \binom{|d|}{\lceil|d|/2\rceil}\frac{1}{2^{|d|}}$$

### A.3  Additional results

In the given section we provide the figures, providing a comprehensive compression analysis comparing TF-IDF and LDA techniques across three distinct datasets (IMDB, AG News, and arXiv). The analysis evaluates three key metrics - Compression Ratio (CR), Jensen-Shannon Divergence (JSD), and ROUGE-F1 scores - as functions of dictionary compression quantile p, with results plotted against their theoretical estimators. The green shaded regions represent confidence intervals around the estimated values, while the black dots indicate the true theoretical values for comparison. Both TF-IDF (top three rows) and LDA (bottom three rows) methods show varying performance patterns across the different datasets, with the estimation curves generally tracking well with their corresponding theoretical benchmarks.

Figure 3: Compression analysis for TF-IDF (top three rows) and LDA (bottom three rows) techniques. The results compare the compression ratio CR, Jensen-Shannon divergence JSD, and ROUGE-F1 scores, as functions of the dictionary compression quantile $p$, with their theoretical estimators across the IMDB, AG News, and arXiv datasets .

## A.4 Experiment algorithms

Here, we describe the algorithms referenced in the main text and used throughout the experimental section. For both of the central components of the paper – the analysis of compression-based representations and the evaluation of statistical bounds – we provide clear pseudo-code that can be directly translated into practical implementations. The goal of presenting the algorithms in the appendix is to give the reader a transparent view of how the theoretical quantities are computed in practice, bridging the gap between abstract definitions and experimental procedures. Each algorithm is written in a way that emphasizes the logical flow of operations, starting from the input dataset, applying compression or transformation, and proceeding to the estimation of key quantities such as divergences, bounds, and error measures. By doing so, we aim to highlight that the computational steps are straightforward and reproducible, and that they can be adapted to other datasets or models with minimal modification.

---

**Algorithm 1** Clusterization statistics collection

---

**Input:** Dataset $X$, compression model $f_{comp} \in \{\text{tf-idf}, \text{LDA}\}$, $p_{values}$ list of possible compression parameters.

**Return:** $D_p$ dictionary of statistics.

$D_p \leftarrow \{\}$

**for** $p$ in $p_{values}$ **do**

    $X_c \leftarrow f_{comp}(X, p)$

    $\widehat{Y}_p \leftarrow Stats(X_c, p)$ {Calculate statistics based on Theorems 1 – 4 with $X_c$}

    $Y_p \leftarrow TrueValues(X_c, p)$ {Calculate true values based on definitions in Section 3.1.}

    $D_p[p] \leftarrow (\widehat{Y}_p, Y_p)$ {Save the bounds and estimators for the given value of $p$}

**end for**

---

---

**Algorithm 2** Encoding statistics collection

---

**Input:** Dataset $X$, dimension size $D$, $epochs$ number of epochs of Monte Carlp, compression model $f_{comp} \in \{\text{tf-idf}, \text{LDA}\}$, $p_{values}$ list of possible compression parameters.

**Return:** $E$ the list of encoding statistics

$E \leftarrow []$

**for** $i$ in $[1, \ldots, epochs]$ **do**

    $\Phi(\mathcal{A}) \leftarrow U(\{\pm 1\}^{|\mathcal{A} \times d|})$ {Generate random vectors}

    $\widehat{\varepsilon}_p \leftarrow \{p : 0\}$ {Dict for interesting values of p}

    **for** $j$ in $[1, \ldots, epochs]$ **do**

        **for** $p$ in $p_{values}$ **do**

            $d_1', d_2' \leftarrow f_{comp}(d_1, p), f_{comp}(d_2, p)$ {Compress the documents}

            $\phi_1', \phi_2' \leftarrow \phi(d_1'), \phi(d_2')$ {Encode the documents}

            $\widehat{\varepsilon}_p[p] = \widehat{\varepsilon}_p[p] + \dfrac{|\langle \phi_1', \phi_2' \rangle|}{D}$

        **end for**

    **end for**

    $\widehat{\varepsilon}_p[p] = \widehat{\varepsilon}_p[p]/epochs$ {Average the value of dot-product}

    $E = E \cup \widehat{\varepsilon}_p$

**end for**

$E = (\text{mean}(E), \text{std}(E))$ {Average and get std of all estimators}

---

# ASR Models for Traditional Emirati Arabic: Challenges, Adaptations, and Performance Evaluation

**Maha AlBlooki**
Mohamed bin Zayed University of AI
Abu Dhabi, UAE
maha.alblooki@mbzuai.ac.ae

**Kentaro Inui**
Mohamed bin Zayed University of AI
Abu Dhabi, UAE
kentaro.inui@mbzuai.ac.ae

**Shady Shehata**
University of Waterloo
Ontario, Canada
shady.shehata@uwaterloo.ca

## Abstract

Traditional Emirati Arabic, a culturally rich and linguistically distinct dialect, remains underrepresented in modern automatic speech recognition (ASR) systems. This paper addresses the gap by introducing a curated speech corpus derived from heritage broadcasts and literary sources, and by evaluating the performance of state-of-the-art ASR models on this low-resource dialect. We examine the zero-shot and fine-tuned performance of five pre-trained models—Wav2Vec2, XLS-R, Whisper, and Massively Multilingual Speech (MMS)—on our traditional Emirati Arabic dataset. Our results show that fine-tuning improves both Word Error Rate (WER) and Character Error Rate (CER), with MMS achieving the best results post-adaptation. Through detailed error analysis, we highlight challenges posed by dialectal morphology, phonology, and lexical variation, and propose targeted adaptations for dialect-specific ASR. This work establishes a foundational benchmark for traditional Emirati ASR and contributes to the broader goal of preserving linguistic heritage through speech technology.

## 1 Introduction

Automatic Speech Recognition (ASR) technologies have achieved remarkable performance in high-resource languages such as English and Mandarin. However, their effectiveness diminishes sharply for low-resource languages and dialects, particularly those with significant phonological and morphological variation. Arabic presents unique challenges in this regard, being a highly diglossic language with numerous regional dialects, many of which are underserved by current ASR systems.

Traditional Emirati Arabic is one such dialect. Rooted in the oral traditions of the United Arab Emirates, it retains linguistic features from Bedouin, coastal, and mountain communities that are increasingly overshadowed by Modern Standard Arabic (MSA) and urban Gulf variants. This

dialect is not only linguistically distinct but also culturally significant, encoding idiomatic expressions, heritage knowledge, and regional identity.

Table 1: Linguistic Features of Traditional Emirati Dialect

| Feature Type | Example from Transcript | Description |
|---|---|---|
| Phonological | جـيـه (Chaih), مب (mub) | ج pronounced as /ch/ (instead of /j/); consonant reduction from ما هو (ma huwa) |
| Morphological | بنتعاون, بنخبر, يسوونه | Prefix ب for future tense; Gulf-specific plural verb conjugation |
| Lexical | اللّوليين, هلنا, ربعنا, يارنا | Heritage terms for "our family", "elders", "neighbors", "our friends" |
| Syntactic | يـوم بتكـبر هالشي | Use of يـوم (yawm) for conditionals; contracted demonstrative هالشي |
| Discourse Markers | يعني, على طول | Filler word يـعـني ; Gulf expression على طـول meaning "immediately" |

Despite its value, traditional Emirati Arabic has been largely ignored in computational linguistics. Existing ASR systems are ill-equipped to handle its unique phonetic and lexical traits. To address this gap, we develop a dedicated speech dataset sourced from the *Alsanaa* (Dalmook, 2021) program and related literary content, and evaluate how modern ASR models perform on this data.

In this paper, we present:

- A curated traditional Emirati speech corpus with standardized transcription and preprocessing

- A comparative evaluation of five leading ASR models (Wav2Vec2, XLS-R, Whisper Small, Whisper Medium, MMS) in zero-shot and fine-tuned scenarios

- Insights into model-specific strengths and limitations for dialectal Arabic ASR.

## 2 Related Work

### 2.1 ASR for Arabic and Dialectal Variants

ASR systems have achieved remarkable progress for major world languages, yet robust solutions for Arabic dialects-particularly traditional Emirati Arabic-remain limited due to unique linguistic features and data scarcity. The Emirati dialect, with its distinct phonological and grammatical characteristics, poses significant challenges for ASR, especially given the lack of dedicated speech resources. Addressing such dialectal diversity is crucial for both technological inclusion and cultural preservation.

Recent advances in self-supervised learning (SSL) have enabled substantial improvements in ASR for low-resource languages and dialects. Models such as wav2vec2, HuBERT, and WavLM have demonstrated strong performance gains when fine-tuned on limited labeled data (Zhao and Zhang, 2022). Cross-lingual models, including XLS-R and Meta's MMS, further extend these capabilities, with XLS-R achieving impressive results even with as little as five minutes of training data in Indonesian language experiments (Sakti and Titalim, 2023). For Arabic, multilingual SSL models generally outperform monolingual approaches, as shown by Younis and Mohammad (2023), who report that fine-tuned XLS-R and MMS models achieve lower word error rates (WER) compared to monolingual baselines.

End-to-end models such as Whisper have also gained traction for their ability to generalize across languages. Talafha et al. (2023) benchmarked Whisper on multiple Arabic dialects, finding that while zero-shot performance often surpasses fully fine-tuned XLS-R models, significant drops occur for previously unseen dialects, including Emirati. The VoxArabica system further demonstrates the potential of SSL-based models for both dialect iden-

tification and ASR across a wide range of Arabic varieties (Waheed et al., 2023).

Hybrid approaches that combine deep learning with traditional phonetic modeling have also been explored. Dhouib et al. (2022) provide a systematic review of Arabic ASR research, highlighting the predominance of MSA-focused studies and the underrepresentation of dialectal variants. Novel architectures, such as CNN-LSTM with attention mechanisms, have shown promise for dialectal ASR, with Alsayadi et al. (2022) reporting improved WER on SASSC and MGB-3 datasets.

### 2.2 Low-Resource ASR Techniques

Transfer learning is a key strategy for improving ASR in low-resource settings. Elmahdy et al. (2014) utuilize MSA data to enhance recognition of under-resourced Arabic dialects, achieving notable WER reductions for Qatari Arabic. Data augmentation methods, including SpecAugment, synthetic speech, and self-training, have also proven effective. Bartelds et al. (2023) demonstrate that self-training and TTS-based augmentation consistently reduce WER for minority languages. Similarly, Khudhair and Talib (2022) show that combining data augmentation with language modeling yields competitive results for Arabic ASR.

Innovative data creation pipelines further address resource scarcity. Yeroyan and Karpov (2024) introduce a workflow for generating ASR datasets from audiobooks, enabling practical ASR development for languages with limited training data.

### 2.3 Datasets and Benchmarking

The development of high-quality datasets is foundational for Arabic ASR research. The Casablanca dataset covers eight Arabic dialects, including Emirati, and provides comprehensive annotations for benchmarking (Talafha et al., 2024). Mixat offers Emirati-English code-switching data, highlighting the challenges of bilingual ASR (Ali and Aldarmaki, 2024). SADA (Alharbi et al., 2024) and QASR (Mubarak et al., 2021) further expand resources for Gulf and multi-dialect Arabic speech, supporting supervised training and a range of speech and NLP tasks.

Efforts to benchmark code-switching ASR are exemplified by Hamed et al. (2022), who introduce a new Egyptian Arabic-English corpus and demonstrate the benefits of combining DNN-hybrid and Transformer approaches. Despite these advances,

challenges remain in achieving consistent evaluation and broad dialectal coverage.

## 2.4 Gaps and Motivation

While recent work has advanced ASR for Arabic and its dialects, systematic evaluation and adaptation of state-of-the-art pre-trained models for traditional Emirati Arabic remain largely unexplored. This study addresses this gap by benchmarking and fine-tuning leading ASR models on Emirati speech, aiming to identify effective strategies for robust dialectal ASR and contribute to the broader field of low-resource speech technology.

## 3 Dataset

To develop and evaluate ASR models for traditional Emirati Arabic, we curated a dialect-specific speech corpus sourced from *Alsanaa* (Dalmook, 2021) program, broadcast by *Aloula* station and supported by the Hamdan bin Mohammed Heritage Center. The dataset includes 102 MP3 audio files (approximately 4 hours) and their corresponding transcriptions, extracted from *alsanaa* book, a heritage literature book authored by Abdullah Bin Dalmook. These recordings capture authentic Emirati Arabic speech, preserving the dialectal nuances and linguistic patterns unique to the region.

Given the lack of existing Emirati ASR corpora, we aligned the audio and text manually, converting them into structured plain-text pairs. The dataset was partitioned into 80% training, 10% validation, and 10% test splits. Notably, the audio is spoken by a single male speaker, limiting speaker diversity but preserving dialectal authenticity.

من الأشياء الجميله إللي شفناها
عند هلنا الأوليين إنهم يصالحون
بين العرب. يعني إذا اثنينه متزاعلين
ولاّ اثنينه متضاريين ولاّ شي،
ساروا وصالحوا بينهم.

Figure 1: Sample transcription

Preprocessing included:

- Diacritics and punctuation removal to standardize transcriptions

- Audio cropping (removal of non-speech intro/outro segments)

- Mono conversion and resampling to 16 kHz

- Normalization to standardize amplitude levels

This dataset captures phonological, morphological, and lexical features unique to traditional Emirati Arabic and serves as a foundational resource for dialect-specific ASR. The full dataset and preprocessing pipeline are available online.[1]

## 4 Models and Training

We adopt a comparative experimental framework to evaluate the performance of state-of-the-art ASR models on traditional Emirati Arabic. Our approach consists of two main stages: zero-shot evaluation and fine-tuning.

### 4.1 Model Selection

We evaluated five pre-trained ASR architectures, each fine-tuned or adapted for Emirati Arabic or closely related dialects:

**Wav2Vec 2.0 (`eabayed/wav2vec2emiratidialict`[1])** Wav2Vec 2.0 is a self-supervised learning framework for speech recognition that learns audio representations via a contrastive task, enabling strong performance with limited labeled data (Baevski et al., 2020). Its architecture combines a convolutional feature encoder with a Transformer network, allowing effective modeling of phonetic and lexical features in low-resource settings. The model used here is further adapted to Emirati Arabic using audio from regional media, resulting in 315 million parameters and improved recognition of dialectal nuances.

**XLS-R (`jonatasgrosman/wav2vec2-large-xlsr-53-arabic`)** XLS-R extends Wav2Vec 2.0 to the multilingual domain, pre-trained on over 436,000 hours of speech in 128 languages (Babu et al., 2021). This enables robust cross-lingual transfer and strong performance on low-resource dialects. The Arabic-adapted variant, with 315 million parameters, is fine-tuned on Common Voice 6.1 and the Arabic Speech Corpus, making it well-suited for Emirati Arabic (Babu et al., 2021).

**Whisper Small (`ayoubkirouane/whisper-small-ar`)** Whisper is a transformer-based encoder-decoder ASR system trained on diverse multilingual data

---

[1] https://github.com/MahaAlBlooki/alsanaa-emirati-dataset

([Radford et al., 2022](#)). The small Arabic model (241M parameters) is fine-tuned on the Mozilla Common Voice v11 dataset for Arabic, and further adapted for Emirati speech, balancing efficiency with accuracy.

**Whisper Medium (`Seyfelislem/whisper-medium-arabic`)** This variant of Whisper, with 763 million parameters, is optimized for Arabic speech recognition. Fine-tuning on Emirati data enhances its ability to transcribe dialectal speech, leveraging the robust encoder-decoder architecture of Whisper.

**MMS (`facebook/mms-1b-all`)** Massively Multilingual Speech (MMS) is a self-supervised model trained on over 1,000 languages, including Arabic dialects ([Zhang et al., 2023](#)). With 965 million parameters, MMS is designed for broad language coverage and demonstrates strong zero-shot and few-shot ASR capabilities. While not specifically fine-tuned for Emirati Arabic, its multilingual training enables generalization to underrepresented dialects.

Each model was evaluated in two modes:

- **Zero-shot inference**: Direct evaluation without further training on our dataset.

- **Fine-tuning**: Models were adapted to the Emirati dataset using transfer learning.

## 4.2 Fine-Tuning Strategy

Fine-tuning involved freezing most pretrained layers and training only the final layers (e.g., projection heads and classification layers). The following configuration was used:

- **Optimizer**: AdamW with weight decay

- **Learning rate schedule**: Linear warm-up followed by decay

- **Batch size**: Adjusted per model based on memory constraints

- **Epochs**: Trained until validation loss convergence (early stopping applied)

- **Data augmentation**: Speed perturbation and SpecAugment to improve generalization

- **Gradient accumulation**: Enabled to simulate larger batch sizes on limited hardware

## 5 Evaluation

### 5.1 Metrics

We use two standard ASR metrics:

- **Word Error Rate (WER)**: Percentage of word-level errors (insertions, deletions, substitutions).

- **Character Error Rate (CER)**: Measures character-level discrepancies; useful for morphologically rich languages and dialects.

Both metrics were calculated on the validation and test splits of our Emirati dataset.

### 5.2 Evaluation Protocol

All models were tested directly on the test set without any adaptation to measure out-of-the-box generalization. After training, models were evaluated on the same test set to assess improvements in recognition accuracy.

### 5.3 Qualitative Analysis

Beyond quantitative metrics, we conducted a qualitative error analysis focused on the recognition of dialect-specific lexical items, morphological transformations (e.g., future tense prefixes), and common phonological shifts (e.g., hamza deletion, /j/ → /ch/ substitutions).

## 6 Results

We evaluated the zero-shot and fine-tuned performance of several state-of-the-art pre-trained ASR models-Wav2Vec 2.0, XLS-R, Whisper (small and medium), and MMS-on traditional Emirati Arabic speech. Performance was measured using WER and CER, providing insight into both word-level and subword recognition accuracy.

### 6.1 Baseline Performance

In the zero-shot setting in Table [2](#), Wav2Vec 2.0 achieved the best results among all models, with a WER of 46.50% and CER of 17.13%. This suggests that its self-supervised pre-training enables effective generalization to unseen dialects, capturing phonetic patterns even when word-level recognition is challenging. MMS ranked second (WER 67.21%, CER 24.56%), likely benefiting from its broad multilingual training and explicit support for Arabic dialects. XLS-R, despite its cross-lingual design, performed poorly (WER 88.26%, CER

| Model | WER (%) | CER (%) |
|---|---|---|
| Wav2Vec 2.0 | 46.50 | 17.13 |
| XLS-R | 88.26 | 40.37 |
| Whisper Small | 93.06 | 81.02 |
| Whisper Medium | 86.10 | 75.01 |
| MMS | 67.21 | 24.56 |

Table 2: Average WER and CER on the whole dataset in baseline inference

40.37%), indicating potential limitations in its coverage of Gulf Arabic and a significant domain gap when applied to Emirati speech. Whisper models showed the weakest zero-shot performance, with Whisper Small reaching 93.06% WER and 81.02% CER, and Whisper Medium slightly better at 86.10% WER and 75.01% CER. The high CER values for Whisper indicate substantial difficulties at the character level, likely due to mismatches between the pre-training data and the phonological characteristics of Emirati Arabic.

Qualitative analysis of model errors revealed that models often misrecognized dialect-specific vocabulary and morphemes, with XLS-R and Whisper in particular producing transcriptions influenced by other Arabic dialects. For example, XLS-R frequently substituted Emirati morphemes with those more typical of Egyptian or Levantine Arabic, reflecting gaps in dialectal representation in the pre-training corpus.

These results highlight the challenges of recognizing traditional Emirati Arabic with existing ASR models and underscore the importance of dialect-specific adaptation. The findings establish a benchmark for future work and inform model selection and adaptation strategies for low-resource dialectal ASR, with broader implications for Arabic speech technology research

### 6.2 Fine-Tuned Performance

Table 3 summarizes the impact of fine-tuning each ASR model on the Emirati *Alsanaa* dataset. Fine-tuning led to substantial performance gains for some architectures, while others showed limited or even negative adaptation.

MMS exhibited the most pronounced improvement, with WER dropping from 67.21% to 41.04% and CER from 24.56% to 13.34%. This 26.17 and 11.22 percentage point reduction in WER and CER, respectively, highlights the effectiveness of MMS's multilingual pre-training in facilitating rapid adaptation to low-resource dialects. After fine-tuning, MMS outperformed all other models, establishing a new benchmark for Emirati Arabic ASR.

Wav2Vec 2.0 also benefited from fine-tuning, achieving a modest reduction in WER (from 46.50% to 44.30%) and CER (from 17.13% to 15.96%). The relatively small improvement suggests that the model's self-supervised representations already captured much of the dialectal variation present in the dataset, resulting in stable performance before and after adaptation.

In contrast, XLS-R's performance deteriorated after fine-tuning, with WER rising from 88.26% to 89.78% and CER from 40.37% to 42.31%. This decline may indicate overfitting to the limited training data or challenges in adapting broad cross-lingual representations to specific dialectal features, a phenomenon noted in low-resource ASR adaptation literature.

The Whisper models showed mixed results. Fine-tuning Whisper Small led to further degradation, with WER exceeding 100% (100.04%); it seemed like the model was encountering a repetition or loop behavior at the end of some transcriptions. For instance, in one of the transcriptions, تومیه, a gibberish prediction of what is supposed to be توايهوا (*cheek-kissed*) is repeated many times consecutively, which isn't in the original text. This type of repetition has artificially inflated the WER of Whisper Small model. On the other hand, the CER increased to 75.53%, suggesting substantial insertion errors and a mismatch between model architecture and the Emirati dialect under data-scarce conditions. Whisper Medium showed only marginal change, with WER shifting from 86.10% to 88.60% and persistently high CER, indicating that additional data or specialized adaptation techniques may be required for effective dialectal ASR with Whisper.

Overall, these results underscore the importance of model selection and adaptation strategy for low-resource dialectal ASR. While MMS demonstrates strong adaptability to Emirati Arabic, other architectures may require more sophisticated fine-tuning or larger datasets to achieve competitive performance.

### 6.3 Error Analysis

A detailed error analysis reveals notable differences in how each model adapts to traditional Emirati Arabic, highlighting both architectural

كالشي يعدنا هنيه له سلام خاصبه

و اليوم انت اندخلت على عرب

و ريته و يتغدون يحتاي

ان تقول لهم هنهم

و هي دعوه من الله ان يهنيههم بالأكل

و هم يحتاي يردون عليهك

و يقولون و منهم و المعنى انه انت

قربوا يا نا و تغدوا يا نا

و تانه هم مُتغدّي

و اللي مبماكلت را بيتغدّه

يا عمو تاني متغدّي ما بيرو مزيد

لكنه الاصل انه هذا هو الشلام

و السلام عقب الغداء انه تغدّه و خلّص

و انه ما تغدّه و خلّص

و انه جافضه من الغداء نشو

توميه و توميه و توميه و توميه

و توميه و توميه و توميه و توميه

و توميه و توميه و توميه و توميه

و توميه و توميه و توميه و توميه

Figure 2: Sample transcription with decoding repetition

strengths and persistent challenges. Models based on self-supervised pre-training, such as MMS and Wav2Vec 2.0, consistently outperformed Whisper variants, suggesting that phonetic representation learning is more effective for dialectal ASR than multitask training approaches.

A key observation is the relationship between zero-shot and fine-tuned performance: models with strong zero-shot results (e.g., Wav2Vec 2.0) exhibited only modest improvements after fine-tuning, while models with moderate zero-shot performance (e.g., MMS) showed substantial gains. This pattern

| Model | WER (%) | CER (%) |
|-------|---------|---------|
| Wav2Vec 2.0 | 44.3 | 15.96 |
| XLS-R | 89.78 | 42.31 |
| Whisper Small | 100.04 | 75.53 |
| Whisper Medium | 88.60 | 72.03 |
| MMS | 41.04 | 13.34 |

Table 3: Average WER and CER on test set after fine-tuning

underscores the importance of evaluating both generalization and adaptability when selecting ASR architectures for low-resource dialects.

Across all models, CER was consistently lower than WER, indicating that character-level recognition is more robust than word-level recognition. This discrepancy, especially pronounced in Wav2Vec 2.0 and MMS, suggests that while phonetic patterns are captured effectively, models struggle with accurate word segmentation and lexical reconstruction. Integrating language models during post-processing may help mitigate these issues.

Architectural differences also affected data efficiency. MMS demonstrated high data efficiency, achieving significant improvements with limited Emirati data, whereas XLS-R and Whisper required more extensive adaptation to yield comparable results. Notably, fine-tuned Whisper Small frequently truncated longer utterances, omitting culturally salient content and narrative details. Additionally, repetition errors were observed, with the model generating nonsensical word sequences, artificially inflating the WER.

Dialectal specialization remains a significant challenge. Even after fine-tuning, high error rates persisted-particularly for Whisper Small and XLS-R, which are primarily pre-trained on Egyptian or MSA data. These models often substituted Emirati morphemes with forms from other dialects, reflecting insufficient representation of Gulf Arabic in the pre-training corpus. Furthermore, inconsistent diacritization in XLS-R outputs, despite ground-truth normalization, introduced additional errors.

These findings emphasize the need for careful model selection, larger dialectal datasets, and potentially pre-training strategies tailored to Gulf Arabic. The persistent performance gaps highlight the ongoing challenge of developing inclusive ASR technologies for underrepresented dialects, underscoring the importance of both technical innovation and investment in dialectal language resources.

## 7 Limitations

This work faces several limitations, one of which is dataset diversity. The dataset includes a single speaker (male), limiting phonetic and demographic diversity. This may bias model performance toward that speaker's vocal and dialectal traits. Another limitation is duration. With only 4 hours of audio, the dataset is relatively small, constraining model generalization. Additionally, dialect cover-

age is limited. While rich in traditional features, the dataset does not fully represent all sub-dialectal varieties across the UAE (e.g., eastern vs. western tribal variants). Moreover, the evaluation scope of WER and CER focused on transcription accuracy, without assessing downstream tasks, such as speaker identification or sentiment analysis.

Future work should explore speaker diversity, cross-dialectal robustness, and larger-scale datasets.

# 8 Conclusion

This paper presents the first ASR benchmark for traditional Emirati Arabic, a linguistically and culturally significant but technologically underserved dialect. By compiling a novel dataset and evaluating state-of-the-art ASR models in both zero-shot and fine-tuned settings, we demonstrate the value of transfer learning and domain-specific adaptation.

Our results demonstrate that self-supervised models with strong multilingual pre-training, particularly MMS, achieve superior adaptability and performance after fine-tuning, while other architectures exhibit varying degrees of success. The persistent gap between character- and word-level accuracy underscores the need for improved modeling of dialectal lexical and phonological features.

This work contributes to Arabic dialectal ASR research and highlights the role of speech technology in preserving oral heritage. We release our dataset and preprocessing tools to encourage further research on Gulf Arabic ASR.

# 9 Ethics Statement

This research adheres to the ACL Ethics Policy. All audio recordings used in this study were publicly available and sourced from cultural heritage broadcasts and literary materials produced by the Hamdan bin Mohammed Heritage Center. Proper credit has been given to the original author, Abdullah Bin Dalmook, whose work was used with the intent of preserving linguistic and cultural heritage.

The dataset features speech from a single speaker who is a public broadcaster and author. No personally identifiable or sensitive information is included. The goal of this research is to support inclusive technology and cultural preservation, not surveillance or misuse.

We acknowledge the potential risks of dialectal ASR systems being misused for sociolinguistic profiling or discrimination. To mitigate this, our work is released with a cultural preservation focus, encouraging ethical use in academic and heritage documentation contexts.

# References

Sadeen Alharbi, Areeb Alowisheq, Zoltán Tüske, Kareem Darwish, Abdullah Alrajeh, Abdulmajeed Alrowithi, Aljawharah Bin Tamran, Asma Ibrahim, Raghad Aloraini, Raneem Alnajim, Ranya Alkahtani, Renad Almuasaad, Sara Alrasheed, Shaykhah Alsubaie, and Yaser Alonaizan. 2024. Sada: Saudi audio dataset for arabic. In *ICASSP 2024 - 2024 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, pages 10286–10290.

Maryam Al Ali and Hanan Aldarmaki. 2024. Mixat: A data set of bilingual emirati-english speech.

Hamzah A. Alsayadi, Salah Al-Hagree, Fahd A. Alqasemi, and Abdelaziz A. Abdelhamid. 2022. Dialectal arabic speech recognition using cnn-lstm based on end-to-end deep learning. In *2022 2nd International Conference on Emerging Smart Technologies and Applications (eSmarTA)*, pages 1–8.

Arun Babu, Changhan Wang, Andros Tjandra, Kushal Lakhotia, Qiantong Xu, Naman Goyal, Kritika Singh, Patrick von Platen, Yatharth Saraf, Juan Pino, Alexei Baevski, Alexis Conneau, and Michael Auli. 2021. XLS-R: self-supervised cross-lingual speech representation learning at scale. *CoRR*, abs/2111.09296.

Alexei Baevski, Henry Zhou, Abdelrahman Mohamed, and Michael Auli. 2020. wav2vec 2.0: A framework for self-supervised learning of speech representations. *CoRR*, abs/2006.11477.

Martijn Bartelds, Nay San, Bradley McDonnell, Dan Jurafsky, and Martijn Wieling. 2023. Making more of little data: Improving low-resource automatic speech recognition using data augmentation.

Abdullah Hamdan Bin Dalmook. 2021. *Alsanaa*. Hamdan bin Mohammad Heritage Center, Dubai. Paper Cover; Dimensions: 24x17 cm.

Amira Dhouib, Achraf Othman, Oussama El Ghoul, Mohamed Koutheair Khribi, and Aisha Al Sinani. 2022. Arabic automatic speech recognition: A systematic literature review. *Applied Sciences*, 12(17).

Mohamed Elmahdy, Mark Hasegawa-Johnson, and Eiman Mustafawi. 2014. Development of a TV broadcasts speech recognition system for qatari Arabic. In *Proceedings of the Ninth International Conference on Language Resources and Evaluation (LREC'14)*, pages 3057–3061, Reykjavik, Iceland. European Language Resources Association (ELRA).

Injy Hamed, Pavel Denisov, Chia-Yu Li, Mohamed Elmahdy, Slim Abdennadher, and Ngoc Thang Vu.

2022. Investigations on speech recognition systems for low-resource dialectal arabic–english code-switching speech. *Computer Speech Language*, 72:101278.

Mohanad Khudhair and Ahmed Talib. 2022. Improving low resources arabic speech recognition using data augmentation. In *2022 Fifth College of Science International Conference of Recent Trends in Information Technology (CSCTIT)*, pages 60–65.

Hamdy Mubarak, Amir Hussein, Shammur Absar Chowdhury, and Ahmed Ali. 2021. Qasr: Qcri al-jazeera speech resource – a large scale annotated arabic speech corpus.

Alec Radford, Jong Wook Kim, Tao Xu, Greg Brockman, Christine McLeavey, and Ilya Sutskever. 2022. Robust speech recognition via large-scale weak supervision.

Sakriani Sakti and Benita Angela Titalim. 2023. Leveraging the multilingual indonesian ethnic languages dataset in self-supervised models for low-resource asr task. In *2023 IEEE Automatic Speech Recognition and Understanding Workshop (ASRU)*, pages 1–8.

Bashar Talafha, Karima Kadaoui, Samar Mohamed Magdy, Mariem Habiboullah, Chafei Mohamed Chafei, Ahmed Oumar El-Shangiti, Hiba Zayed, Mohamedou cheikh tourad, Rahaf Alhamouri, Rwaa Assi, Aisha Alraeesi, Hour Mohamed, Fakhraddin Alwajih, Abdelrahman Mohamed, Abdellah El Mekki, El Moatez Billah Nagoudi, Benelhadj Djelloul Mama Saadia, Hamzah A. Alsayadi, Walid Al-Dhabyani, Sara Shatnawi, Yasir Ech-Chammakhy, Amal Makouar, Yousra Berrachedi, Mustafa Jarrar, Shady Shehata, Ismail Berrada, and Muhammad Abdul-Mageed. 2024. Casablanca: Data and models for multidialectal arabic speech recognition.

Bashar Talafha, Abdul Waheed, and Muhammad Abdul-Mageed. 2023. N-shot benchmarking of whisper on diverse arabic speech recognition.

Abdul Waheed, Bashar Talafha, Peter Sullivan, Abdel-Rahim Elmadany, and Muhammad Abdul-Mageed. 2023. Voxarabica: A robust dialect-aware arabic speech recognition system.

Ara Yeroyan and Nikolay Karpov. 2024. Enabling asr for low-resource languages: A comprehensive dataset creation approach.

Hiba Adreese Younis and Yusra Faisal Mohammad. 2023. Arabic speech recognition based on self supervised learning. In *2023 16th International Conference on Developments in eSystems Engineering (DeSE)*, pages 528–533.

Yunbo Zhang, Deepak Gopinath, Yuting Ye, Jessica Hodgins, Greg Turk, Jungdam Won, Harrison Jesse Smith, Qingyuan Zheng, Yifei Li, Somya Jain, Jessica K. Hodgins, Simran Arora, Patrick Lewis, Angela Fan, Jacob Kahn, Christopher Ré, and Michael Auli. 2023. Scaling speech technology to 1,000+ languages. *arXiv preprint arXiv:2305.13516*.

Jing Zhao and Wei-Qiang Zhang. 2022. Improving automatic speech recognition performance for low-resource languages with self-supervised models. *IEEE Journal of Selected Topics in Signal Processing*, 16(6):1227–1241.

# Style-Controlled Response Generation for Dialog Systems with Intimacy Interpretation

**Takuto Miura, Kiyoaki Shirai, Natthawut Kertkeidkachorn**
Japan Advanced Institute of Science and Technology, Nomi, Ishikawa, Japan
{s2460005, kshirai, natt}@jaist.ac.jp

## Abstract

This paper proposes a novel method to control the style of the dialog system's utterances according to the user's level of intimacy with the system. Specifically, the dialog model generates responses in a polite style when the user exhibits a low level of intimacy with the system and in a casual style when the user's intimacy is high. The proposed model consists of two submodels: the Intimacy Interpreter and Response Generator. The Intimacy Interpreter generates an embedding that represents the user's intimacy. This model is trained by contrastive learning using an intimacy-labeled dialog corpus. The Response Generator accepts a dialog context and an intimacy embedding, and then generates a response in an appropriate style. We apply two loss functions to fine-tune a Large Language Model (LLM) to train the Response Generator. The results of automatic and human evaluations show that the proposed method outperforms the baselines in terms of style control in response generation.

## 1 Introduction

In recent years, free dialog systems that allow users to converse about any topic have attracted considerable attention (Khatri et al., 2018; Higashinaka et al., 2021; Dinan et al., 2020). These systems need to have a comfortable conversation with the user and establish a long-term friendly relationship to facilitate conversation between the user and the dialog system (Ram et al., 2018).

To establish friendly relationships, humans change their speech style based on their level of intimacy and social connections with others to facilitate smooth communication (Wardhaugh and Fuller, 2021; Hovy, 1987; Silverstein, 2003). This ability is referred to as "style control" hereafter. The style control should also be considered in conversations between a human and a system (Kageyama et al., 2018). Consequently, a free dialog system is required to have the capability for style control.

The goal of this research is to develop a dialog system that dynamically adjusts styles according to the user's feelings toward the dialog system. A typical example of style control is that a speaker uses formal/polite expressions or informal/casual expressions by the relationship with their partner (Aapakallio, 2021; Liu and Kobayashi, 2022). Miura et al. (2024a) reported that speakers tend to use a polite style when intimacy with a partner is low and a casual style when intimacy is high. Therefore, we aim to dynamically recognize the user's level of intimacy through their dialog history and enable the dialog system to flexibly use a polite or casual style when intimacy is low or high.

This paper proposes a model that accurately identifies the user's level of intimacy with the dialog system and generates responses in an appropriate style. An intimacy interpreter is introduced to obtain a user embedding that represents the user's intimacy, and then this embedding is fed into a response generator, which is obtained by fine-tuning a Large Language Model (LLM), as a soft prompt. It enables the dialog system to appropriately control polite and casual styles.

The contributions of this paper are summarized as follows.

- We develop a dialog system that dynamically captures the user's intimacy and adjusts responses to be either polite or casual style accordingly.

- We propose a new framework to obtain an abstract representation of the user's intimacy and incorporate it into a dialog model for style control.

- The effectiveness of the proposed method is demonstrated through automatic and human evaluations.

## 2 Related Work

Methods for generating responses in a particular style have been actively studied. Niu and Bansal (2018) defined such a task, created a model for identifying a speech style, and proposed a method for generating responses in a given style (e.g., a polite or casual style). Gao et al. (2019) proposed a model that generated responses in a given style while maintaining consistency with the dialog context by sharing the latent space between conversational modeling and style modeling. Zhu et al. (2021) assumed that conversational modeling and style modeling are contradictory, and proposed a method to separate the representations of content and style within the shared latent space proposed by Gao et al. (2019), where each is represented in different dimensions of the latent space. Zheng et al. (2021) proposed a method for automatically constructing a dialog corpus containing utterances in a given style, which was used to train a dialog model that generated responses in line with the specified style. Specifically, they trained a Seq2Seq (Sequence-to-Sequence) model that transformed a sentence into an equivalent sentence in the specified style using a text corpus of that style. A new dialog corpus was constructed by converting the style of utterances in an original dialog corpus using the trained style conversion model. Yang et al. (2020) proposed STYLEDGPT to fine-tune a pretrained language model to obtain a dialog model that generates utterances in the target style. They designed loss functions for fine-tuning, which were based on a language model of a given style and a classification model for identifying the style of an utterance.

In recent years, several studies have leveraged the text generation capabilities of rapidly advancing LLMs to address style control. Konen et al. (2024) controlled a style in text generation by adding style vectors to the activation of hidden layers in an LLM. Two types of style vectors were proposed: the training-based and activation-based style vectors. The former trained the style vectors using the cross-entropy loss between the output of the LLM for the empty input token and the target sentence. The latter employed the activation vectors of the layers in the LLM for the given target sentences to obtain the style vector. Li et al. (2024) created a dialog corpus containing utterances in 38 different style categories using an LLM, allowing fine-grained styles to be handled in dialog system development.

First, a prompt including the name of a target style is given to the LLM to generate a description of the style and an example sentence. Next, the style description and the example sentence were given to the LLM to generate a rationale that the style of the sentence was consistent with the given style description. Finally, the style name, style description, example sentences, and style rationale as well as a plain context were provided to the LLM to generate a response to the given context in the target style. The constructed dialog corpus consisted of the pairs of the input contexts and the generated responses in different styles.

Although the aforementioned studies can generate natural responses in a specific style, they are limited to considering a single style in style control. In contrast, this study aims to dynamically control multiple styles based on the user's state.

Miura et al. (2024b) proposed a dialog system that flexibly switched between two different styles, the polite style and the casual style, according to the changes in the user's intimacy with the dialog system. The dialog model was trained to generate responses in the polite style when the user's intimacy is low and in the casual style when the intimacy is high, by referring to the intimacy estimation model and two language models of the polite and casual styles. In addition, the style discrimination model was employed to train a dialog model so that the probability of the polite (or casual) style of generated responses, which was estimated by the style discrimination model, became high when the user's intimacy was low (or high). This learning method succeeded in achieving better style control capability than general dialog models. However, there is much room to improve the accuracy of style control due to the poor performance of the intimacy estimation model incorporated in the dialog model. Therefore, this study aims to develop a model for interpreting the user's intimacy by creating user embeddings, so the model could accurately capture the user's intimacy and appropriately perform style control.

## 3 Proposed Method

### 3.1 Overview

Figure 1 shows an overview of the dialog model that changes the style based on the recognized user's intimacy. Given a dialog history $X$, the proposed system generates $Y$ which is a response to $X$. Here, $X$ is a conversation between a system $S$ and

$Y = S_4$ | That's good to hear.

Response Generator (LLM)

soft prompt + token sequence

$V'$

FCL

$V$

Intimacy Interpreter

$X^u = \{U_1, U_2, U_3\}$

$S_1$: It's been a while.
$U_1$: Yes, it's been a month.
$S_2$: How is work going these days?
$U_2$: Not bad.
$S_3$: Have you gone on that trip?
$U_3$: Yes, I did and it was a lot of fun.

$X = \{S_1, U_1, S_2, U_2, S_3, U_3\}$

Figure 1: Overview of proposed method

a user $U$, denoted as $X = \{S_1, U_1, \cdots, S_n, U_n\}$, while $Y$ is the next utterance of the system, i.e., $Y = S_{n+1}$.

The proposed system consists of two submodels. The first is the Intimacy Interpreter. It takes the user's past utterances $X^u = \{U_1, \cdots, U_n\}$ as input and interprets the user's degree of intimacy with the dialog system. The output of the Intimacy Interpreter is an intimacy embedding, a vector representation of the user's intimacy. The second is the Response Generator, which is based on an LLM. It takes the dialog history $X = \{S_1, U_1, \cdots, S_n, U_n\}$ as input and produces a response $Y$ as output. At the beginning of the input token sequence, a soft prompt of the user's intimacy is added. This is a single token embedding derived from the intimacy embedding. Specifically, the size of the intimacy embedding produced by the Intimacy Interpreter is changed to that of the token embeddings of the LLM by the Fully Connected Layer (FCL). It is expected that the response is generated in a casual style when the user's intimacy is high and in a polite style when it is low. The length of the dialog history is 3 in Figure 1, but it can be changed arbitrarily.

The following sections describe the details of the Intimacy Interpreter and Response Generator, respectively.

## 3.2 Intimacy Interpreter

The Intimacy Interpreter aims to capture the complex and vague nature of the user's intimacy by representing it as an abstract vector. Hereafter, the Intimacy Interpreter is denoted as $P_{II}(V|X^u)$. The

model takes as input the $n$ consecutive utterances of a user in a dialog context, $X^u = \{U_1, \cdots, U_n\}$, and outputs a vector $V$ representing the user's intimacy with the dialog system.

This study applies contrastive learning to train the Intimacy Interpreter. An intimacy-labeled dialog corpus $D_{in}$, where each dialog is labeled with a 5-point Likert scale indicating the level of intimacy of a speaker with a dialog partner, is used for contrastive learning. The details of this corpus are described in 4.1.1. The user's $n$ consecutive utterances in $D_{in}$ are extracted as a sample $(X_i^u, IL_i)$, where $IL_i$ denotes the five-scale intimacy label assigned to the sample $X_i^u$. Two samples $X_i^u$ and $X_j^u$ are randomly taken from the training data. If the intimacy labels $IL_i$ and $IL_j$ assigned to these two samples are the same, the parameters of the Intimacy Interpreter are updated so that the embedded vectors $V_i$ and $V_j$ become similar. If $IL_i$ and $IL_j$ are not equal, the parameters are updated so that $V_i$ and $V_j$ are different. Specifically, the contrastive loss for training $P_{II}(V|X^u)$ is defined as Equation (1).

$$L_I = \begin{cases} 1 - sim_{cos}(V_i, V_j) & \text{if } IL_i = IL_j \\ |IL_i - IL_j| \cdot \max(0, sim_{cos}(V_i, V_j)) \\ & \text{if } IL_i \neq IL_j \end{cases} \quad (1)$$

$sim_{cos}(\cdot, \cdot)$ represents the cosine similarity between the two sample embedding vectors. When $IL_i \neq IL_j$, the loss becomes large when the difference between $IL_i$ and $IL_j$ is large by giving $|IL_i - IL_j|$ as the weight. The Intimacy Interpreter is obtained by fine-tuning the pre-trained BERT (Devlin et al., 2019) using this loss.

## 3.3 Response Generator

The Response Generator is denoted as $P_{RG}(Y|V', X)$, where $X$ is the dialog history, $V'$ is the soft prompt derived from the intimacy embedding ($V$), and $Y$ is the response to be generated. This subsection describes the details of training the Response Generator.

### 3.3.1 Loss for Style Control

As described earlier, the Response Generator is obtained by fine-tuning an LLM. Following the study of (Miura et al., 2024b), two loss functions, the intimacy-aware word-level loss and the intimacy-aware sentence-level loss, are used to fine-tune the LLM so that the Response Generator generates responses in the appropriate style

(polite or casual) according to the user's intimacy. **Preliminary** The intimacy-labeled dialog corpus $D_{in}$ described in subsection 3.2 is also used to train the the Response Generator. In addition, two style corpora are prepared to handle polite and casual styles in response generation. One is $C_{po}$ which consists of polite style sentences, and the other is $C_{ca}$ which consists of the casual sentences.

Before the training of the Response Generator, an intimacy estimation model $P(I|X^u)$ is trained in advance. This model predicts $I$, the user's level of intimacy with a dialog system, given the user's past $n$ utterances ($X^u$) as input. In our model, $I$ is defined as either "low" or "high". The intimacy estimation model is pre-trained using $D_{in}$. Note that this is a different model from the Intimacy Interpreter $P_{II}(V|X^u)$. The Intimacy Interpreter produces the intimacy embedding, while the intimacy estimation model is a binary classifier. **Intimacy-aware Word-Level Loss** Two style language models are pre-trained. A polite style language model $P_{po}(T)$ is trained using $C_{po}$, and a casual style language model $P_{ca}(T)$ is trained using $C_{ca}$. These models evaluate how likely the given sentence $T$ is in the polite or casual style. They are employed to calculate the polite style word-level loss $L_w^{po}$ and the casual style word-level loss $L_w^{ca}$, respectively, as shown in Equation (2).

$$L_w^s = d(\mathbf{p_Y}||\hat{\mathbf{p}}_\mathbf{Y}) \stackrel{\text{def}}{=} \sum_{i=1}^{m} D_{KL}(p_{y_i}||\hat{p}_{y_i}), \quad (2)$$

where $s$ denotes the style, either $po$ (polite) or $ca$ (casual). This loss is computed for each dialog sample $(X, Y)$ in the training data. $Y$ is denoted as a token sequence $\{y_1, \cdots, y_m\}$. Let $\mathbf{p_Y} = \{p_{y_1}, \cdots, p_{y_m}\}$ be the distribution of the predicted probability of the next word given by the dialog model $P_{RG}(Y|V', X)$, and $\hat{\mathbf{p}}_\mathbf{Y} = \{\hat{p}_{y_1}, \cdots, \hat{p}_{y_m}\}$ be the probability distribution predicted by the style language model $P_s(T)$. $D_{KL}$ is the Kullback-Leibler divergence of the two probability distributions, indicating whether the words generated by the dialog model follow the specified (polite or casual) style.

As shown in Equation (3), the intimacy-aware word-level loss is defined as the weighted sum of two losses, where $p(I{=}\text{low}|X^u)$ is the weight for $L_w^{po}$ and $p(I{=}\text{high}|X^u)$ is the weight for $L_w^{ca}$. $p(I{=}\text{low}|X^u)$ and $p(I{=}\text{high}|X^u)$ are the probabilities of the low intimacy and high intimacy classes, respectively, predicted by the intimacy estimation model.

$$L_w^{in} \stackrel{\text{def}}{=} p(I{=}\text{low}|X^u) \cdot L_w^{po} + p(I{=}\text{high}|X^u) \cdot L_w^{ca} \quad (3)$$

It is expected that this loss will cause the Response Generator to generate more polite style tokens when the intimacy is low, and more casual style tokens when the intimacy is high. **Intimacy-aware Sentence-Level Loss** First, we train a style discrimination model $P'(S|T)$ that classifies the style $S$ of a sentence $T$. The style $S$ is either polite or casual. The style discrimination model is pre-trained from training data in which utterances in $C_{po}$ are samples of the polite class and utterances in $C_{ca}$ are samples of the casual class.

Let $\hat{Y}$ be the response generated by $P_{RG}(Y|V', X)$. The style of $\hat{Y}$ is identified using the style discrimination model $P'(S|T)$, and the $p(S{=}\text{polite}|\hat{Y})$ and $p(S{=}\text{casual}|\hat{Y})$, the predicted probabilities of the polite and casual classes respectively, are calculated. The intimacy-aware sentence-level loss $L_s^{in}$ is defined as the weighted sum of the logarithms of these probabilities, as shown in Equation (4).

$$L_s^{in} \stackrel{\text{def}}{=} -p(I{=}\text{low}|X^u) \cdot \log p(S{=}\text{polite}|\hat{Y})$$
$$-p(I{=}\text{high}|X^u) \cdot \log p(S{=}\text{casual}|\hat{Y}) \quad (4)$$

This loss will contribute to making the Response Generator to generate polite (or casual) style sentences when intimacy is low (or high).

### 3.3.2 Negative Log-likelihood Loss

The two losses described in 3.3.1 are designed to maintain style consistency. A model fine-tuned solely by these losses may exhibit inconsistency between the dialog context and the generated response. Therefore, a common loss for training dialog models, the negative log-likelihood loss defined as shown in Equation (5), is also used. The value $p(Y|V', X)$ denotes the probability of the ground-truth response Y in the training data being generated by the Response Generator for a given soft prompt of user's intimacy $V'$ and the dialog context $X$.

$$L_{NLL} = -\log p(Y|V', X) \quad (5)$$

### 3.3.3 Training Objective

The loss for training the Response Generator, $L_D$, is a weighted sum of two losses for style control ($L_w^{in}$ and $L_s^{in}$) and a loss for content generation ($L_{NLL}$) as follows:

$$L_D = \beta_w \cdot L_w^{in} + \beta_s \cdot L_s^{in} + \beta_{NLL} \cdot L_{NLL} \quad (6)$$

The weights $\beta_w$, $\beta_s$, and $\beta_{NLL}$ are hyperparameters.

### 3.4 Training Details

Our entire dialog model, shown in Figure 1, is trained based on two losses: $L_I$ and $L_D$. On the one hand, the parameters of the Intimacy Interpreter $P_{II}(V|X^u)$ are updated using $L_I$. On the other hand, the parameters of the Response Generator $P_{RG}(Y|V', X)$ and the FCL that transforms the dimension of the intimacy embedding are updated using $L_D$.[1] The Response Generator is based on the LLM, which is computationally expensive to fine-tune. Therefore, LoRA (Hu et al., 2022) is applied to fine-tune the Response Generator.

## 4 Experiments

### 4.1 Datasets

#### 4.1.1 Dialog Corpus with Intimacy Label

The JID corpus (Miura et al., 2024a) is used as the intimacy-labeled corpus $D_{in}$. This corpus consists of recorded and transcribed conversations of about 10 minutes between two speakers. For each conversation, the intimacy labels of each of the two speakers are annotated using a five-point Likert scale. The number of subjects who participated in the dialogs is 19, the number of dialogs is 54, and the total number of utterances is 6,984.

The 54 dialogs in the JID corpus are divided into three subsets: a training set of 33 dialogs, a validation set of 9, and a test set of 12. As mentioned in section 3, the dialog model accepts the preceding dialog context of the user and the system, $X = \{S_1, U_1, \cdots, S_n, U_n\}$, as input and generates the subsequent response $S_{n+1}$ as output. Hereafter, the pair of a dialog context and its corresponding response, denoted by $(X, S_{n+1})$, will be referred to as "response instance." The first $n \times 2$ utterances and the next utterance in a dialog are extracted as $(X, S_{n+1})$. One speaker in the corpus is designated as the system and the other as the user. This procedure is then repeated with the utterance shifted one by one to obtain multiple response instances. In this experiment, the context length is set to $n = 3$. The statistics of the dataset are shown in Table 1.

#### 4.1.2 Style Corpus

Two style corpora of the polite and casual style, $C_{po}$ and $C_{ca}$, are required to train style language

---

[1]The blue modules in Figure 1 indicate the models trained with the loss $L_D$.

|  | Training | Validation | Test |
|---|---|---|---|
| Dialog | 33 | 9 | 12 |
| Response Instance | 4,032 | 921 | 1,284 |

Table 1: Statistics of Dataset



次の対話文脈に対して、あなたはBとして応答を生成してください。
*(For the following dialog context, generate a response as B.)*

【Dialog Context】

Figure 2: Template of Zero-shot Prompt



次の対話文脈に対して、あなたはBとして応答を生成してください。
ただし、AがBに抱く親密度を推測して、親密度が低い場合は丁寧なスタイルで、親密度が高い場合はカジュアルなスタイルで応答を生成してください。
*(For the following dialog context, generate a response as B. Guess the level of intimacy A has with B and generate a response in a polite style if the level of intimacy is low and in a casual style if the level of intimacy is high.)*

【Dialog Context】

Figure 3: Template of Style Control Prompt



---------1st step
この対話からAがBに抱く親密度は
*(From this dialog, the level of intimacy that A feels towards B is)*

【Dialog Context】

---------2nd step
次の対話文脈に対して、あなたはBとして応答を生成してください。
ただし、「【output of the first step】」という解釈を踏まえて、親密度が低い場合は丁寧なスタイルで、親密度が高い場合はカジュアルなスタイルで応答を生成してください。
*(With the interpretation of 【output of the first step】, generate responses in a polite style if the level of intimacy is low, and in a casual style if the level of intimacy is high.)*

【Dialog Context】

Figure 4: Template of Two-step Prompt

models and a style discrimination model. The KeiCO corpus (Liu and Kobayashi, 2022) is used as $C_{po}$. This corpus contains utterances using various types of honorific expressions in Japanese. Besides, $C_{ca}$ is constructed by extracting utterances from conversations between speakers who know each other in the BTSJ Japanese Natural Conversation corpus (Usami, 2021). $C_{po}$ and $C_{ca}$ contain 7,324 and 13,521 utterances, respectively.

### 4.2 Experimental Settings

The following methods are compared in the experiment.

- **Zero-shot prompt (Zero-shot)** This method uses an LLM as a dialog model without fine-tuning or prompting for style control. We only give an instruction for generating responses to the input dialog context. The details of the prompt are shown in Figure 2.

- **Zero-shot prompt for style control (Style**

**control prompt)** This method uses a pre-trained LLM as a dialog model, where a prompt is given to instruct the LLM to generate utterances taking the style control into account. The details of the prompt are shown in Figure 3.

- **Two-step prompt (Two-step)** This method uses a pre-trained LLM as a dialog model using two sequential prompts. We first instruct the LLM to infer the user's level of intimacy, and then to generate the system's response in a polite or casual style according to the inferred level of intimacy. See Figure 4 for details.

- **STYLEDGPT** This is a model where the style is controlled by STYLEDGPT (Yang et al., 2020). Specifically, we fine-tune the LLM to generate utterances that are consistent with the style of the entire JID corpus. The style language model is trained on training data from the JID corpus. The style discrimination model, which distinguishes whether an utterance is in the style of the JID corpus, is trained using utterances from the JID corpus as positive samples and sentences from Japanese Wikipedia as negative samples.

- **Ours_auto** This is our proposed method described in section 3.

- **Ours_gold** Our proposed method where the gold intimacy labels in the JID corpus are used instead of the prediction by the intimacy estimation model. When calculating the losses in Equation (3) and (4), $p(I\text{=low}|X^u)$ and $p(I\text{=high}|X^u)$ given as follows.

$$p(I\text{=low}|X^u) = 1 - \frac{IL}{5} \qquad (7)$$
$$p(I\text{=high}|X^u) = \frac{IL}{5} \qquad (8)$$

$IL$ represents the five-level intimacy label assigned to $X^u$ in the JID corpus. This model evaluates our approach of considering the user's intimacy for the appropriate style control under the ideal condition where the user's intimacy is correctly predicted.

### 4.3 Implementation Details

#### 4.3.1 Intimacy Interpreter and Response Generator

The Intimacy Interpreter described in subsection 3.2 is obtained by contrastive learning based on the Japanese BERT model[2], which was pre-trained on large-scale corpora of Japanese Wikipedia and Japanese CC-100.

The Response Generator described in subsection 3.3 is obtained by fine-tuning llm-3-3.7b[3], which is an LLM based on Transformer (Vaswani et al., 2017) and has been trained on various large Japanese datasets. We also adopted llm-3-3.7b as the LLM for other baseline dialog models.

For the hyperparameters during training, the learning rate for the Intimacy Interpreter is $1\text{e}^{-6}$, while that for the Response Generator is $1\text{e}^{-20}$. For both models, the batch size is 4 and the number of epochs is 5. These values were optimized on the validation set according to the StyCor criteria, which will be defined in subsection 4.5. The Adam optimizer was used to learn the models.

The hyperparameters $\beta_w$, $\beta_s$, and $\beta_{NLL}$ in Equation (6) are set to 0.5, 1, and 0.005, respectively. These values are determined so that the influence of the three types of losses is uniform. Specifically, we calculate the average of the absolute value of each of the three losses in the training data and then determine the weight of each loss as the approximate inverse ratio of the average to the minimum value.

### 4.4 Other Submodels

Several submodels are pre-trained before training of the Intimacy Interpreter and Response Generator.

The style language models $P_{po}(T)$ and $P_{ca}(T)$ are obtained by fine-tuning GPT-2. We use the pre-trained model japanese-gpt2-medium[4], which has been trained on a large Japanese dialog dataset. All utterances in $C_{po}$ and $C_{ca}$ are used to train $P_{po}(T)$ and $P_{ca}(T)$, respectively. The learning rate is set to $5\text{e}^{-4}$, the batch size to 4, and the epoch to 20. The Adam optimizer is used to fine-tune the models.

The style discrimination model $P'(S|T)$ is obtained by fine-tuning the Japanese BERT model[2]. A total of 20,575 utterances are used, comprising 7,274 polite utterances in $C_{po}$ and 13,301 casual utterances in $C_{ca}$. The learning rate is set to $1\text{e}^{-7}$, the batch size to 128, and the epoch to 10. The Adam optimizer is used to fine-tune the model. The accuracy of the style discrimination model was 99% when it was evaluated on the 100 test utterances

---

[2]https://huggingface.co/tohoku-nlp/bert-large-japanese-v2

[3]https://huggingface.co/llm-jp/llm-jp-3-3.7b

[4]https://huggingface.co/rinna/japanese-gpt2-medium

(50 polite and 50 casual) that were not used for training.

The intimacy estimation model $P(I|X^u)$ is based on the Japanese BERT model[2]. The JID corpus is used for fine-tuning the BERT. The learning rate is set to $5e^{-6}$, the batch size to 1, and the epoch to 10. The Adam optimizer is used to train the model. The accuracy of the intimacy estimation model on the test data was 69%.

### 4.5 Evaluation Criteria

Both automatic and human evaluations are carried out to assess responses generated by various methods.

#### 4.5.1 Automatic Evaluation

In the automatic evaluation, the quality of the generated responses is evaluated from three perspectives: relevance, diversity, and style. The relevance is measured by BLEU (Papineni et al., 2002) and ROUGE (Lin, 2004). Specifically, the similarity between a generated response and a ground-truth response is evaluated using BLEU-1, BLEU-2, ROUGE-1, and ROUGE-2. The diversity is measured by Distinct-1 (Dist-1) and Distinct-2 (Dist-2), following the experiment of (Li et al., 2016). The style is evaluated by measuring "Style Correlation" (StyCor). The StyCor metric is defined as the correlation between the probability of the casual style $p(S=\text{casual}|Y)$ and the ground-truth level of the intimacy.[5] This correlation is high when both the predicted probability of the casual style and the level of intimacy are high, or both are low (i.e., the probability of the polite style is high and the level of intimacy is low).

#### 4.5.2 Human Evaluation

The quality of the generated responses is evaluated by humans. To reduce the burden on evaluators, STYLEDGPT and Ours_auto are excluded from the human evaluation. A hundred response instances are randomly taken from the test set of the JID corpus. The dialog context $X$ of each response instance is used as input, and a response is generated using the dialog models. Subjects evaluate these responses from the following three perspectives.

- Style Control: Does the response align with the appropriate style for the relationship between the two speakers? Annotators are also

---

[5]The five-scale score is normalized to values between 0 and 1.

instructed to read the dialog context and guess the relationship between the speakers.

- Relevance: Is the content of the response relevant and consistent with the context?

- Fluency: Is the response natural, fluent, and free of grammatical errors?

For each item, the quality of the responses was assessed by giving a score on a 5-point Likert scale from 1 (inappropriate) to 5 (appropriate). Five native Japanese speakers participated in the manual evaluation. Agreement between annotators' scores was measured using Fleiss's kappa (Fleiss and Jacob, 1973).

## 5 Results

### 5.1 Results of Automatic Evaluation

The results of the automatic evaluation are shown in Table 2. The StyCor of Ours_auto and Ours_gold were 0.239 and 0.250, respectively, outperforming other baseline methods. This confirms that the proposed method, which adjusts the style based on the level of intimacy, can effectively control the polite and casual styles. The decrease of StyCor of Ours_auto compared to Ours_gold may be due to the low accuracy of the intimacy estimation model.

In the evaluation of the relevance, STYLEDGPT and our proposed models achieved better BLEU and ROUGE scores than other baselines, since these models are fine-tuned using the JID corpus, which was the same domain as the test data. However, our models performed slightly worse than STYLEDGPT. On the other hand, the diversity (Dist-1 and Dist-2) of all models was high.

Although the BLEU and ROUGE of our method are worse than those of STYLEDGPT, we think that these indicators are only for reference in automatic evaluation. BLEU and ROUGE only evaluate the similarity between the generated and ground-truth responses, while there could be other appropriate responses that are not included in the dataset. On the other hand, our proposed method clearly outperforms STYLEDGPT in terms of StyCor, indicating superior style control capabilities.

To sum up, our models can improve the ability of the style control with a little decrease in relevance.

### 5.2 Results of Human Evaluation

The results of the human evaluation are shown in Table 3. The "Score" column shows the average

| Method | Relevance | | | | Diversity | | Style |
|---|---|---|---|---|---|---|---|
| | BLEU-1 | BLEU-2 | ROUGE-1 | ROUGE-2 | Dist-1 | Dist-2 | StyCor |
| Zero-shot | 0.0483 | 0.0034 | 0.0780 | 0.0044 | 0.942 | 0.978 | 0.164 |
| Style control prompt | 0.0578 | 0.0053 | 0.1014 | 0.0073 | **0.965** | **0.991** | 0.207 |
| Two-step | 0.0575 | 0.0028 | 0.0932 | 0.0041 | 0.946 | 0.984 | 0.162 |
| STYLEDGPT | 0.2520 | **0.1571** | **0.3392** | **0.2108** | 0.925 | 0.935 | 0.171 |
| Ours$_{auto}$ | 0.2067 | 0.1205 | 0.2986 | 0.1725 | 0.895 | 0.900 | 0.239 |
| Ours$_{gold}$ | **0.2544** | 0.1463 | 0.3390 | 0.1999 | 0.925 | 0.930 | **0.250** |

Table 2: Results of Automatic Evaluation

| Method | Style Control | | | Relevance | | | Fluency | | |
|---|---|---|---|---|---|---|---|---|---|
| | Score | $\kappa$ | $p$ | Score | $\kappa$ | $p$ | Score | $\kappa$ | $p$ |
| Zero-shot | 4.31 | 0.50 | $1e^{-6}*$ | 4.13 | 0.50 | 0.086 | 4.51 | 0.68 | $9e^{-11}*$ |
| Style control prompt | 4.34 | 0.51 | $9e^{-5}*$ | 4.22 | 0.53 | 0.553 | 4.63 | 0.72 | $8e^{-7}*$ |
| Two-step | 4.33 | 0.51 | $6e^{-5}*$ | 4.02 | 0.44 | 0.002* | 4.49 | 0.69 | $2e^{-12}*$ |
| Ours$_{gold}$ | 4.61 | 0.60 | – | 4.26 | 0.54 | – | 4.86 | 0.84 | – |

Table 3: Results of Human Evaluation. * means $p < 0.05$.

score of the five subjects, while the "$\kappa$" column indicates Fleiss's $\kappa$. Welch's test is performed to verify whether there was a significant difference in the scores between Ours$_{gold}$ and other methods. The "$p$" column represents the $p$-value of this statistical test.

For Style Control, Ours$_{gold}$ received the highest score. Additionally, significant differences with all other methods were confirmed. This demonstrates the effectiveness of the approach proposed in this study, which considers the user's level of intimacy for the appropriate selection of polite and casual styles. The $\kappa$ value was 0.60, which indicated moderate agreement.

In terms of Relevance, Ours$_{gold}$ achieved the highest score. However, significant differences were only observed when compared to Two-step. The proposed method performed comparably to the baseline methods in generating responses relevant to the dialog context.

The Fluency score of the proposed method was significantly higher than the other models, indicating its superior ability to generate natural utterances.

## 6 Ablation Study

Table 4 shows the results of the ablation study. The Ours-SCL is the model where two intimacy-aware style control losses, $L_w^{in}$ and $L_s^{in}$, are removed from Equation (6). The Ours-II indicates the removal of the Intimacy Interpreter, which is almost equiva-

lent to the dialog model presented in (Miura et al., 2024b).[6] This model is trained using the gold intimacy labels to calculate the loss $L_D$, so the above two models are compared to Ours$_{gold}$.

The results demonstrated that both the use of the style control losses and the incorporation of the Intimacy Interpreter could improve the StyCor score. Especially, a significant decrease was found in Ours-SCL, indicating that the intimacy-aware style control losses are effective in changing the style appropriately. On the other hand, the contribution of the Intimacy Interpreter was rather limited. It should be noted that both the style control losses and the Intimacy Interpreter could also improve the relevance and diversity of the generated responses.

## 7 Conclusion

In this paper, we proposed the novel method to control the style of a dialog system based on the user's level of intimacy. The model that interpreted the user's level of intimacy was incorporated into the dialog model. This Intimacy Interpreter was trained by contrastive learning using the dialog corpus annotated with the intimacy labels. Furthermore, based on the LLM, which had an excellent capability to generate general responses, we applied two loss functions to improve the model's ability to control the style. The results of both au-

---

[6]The base LLMs are different: llm-jp-3-3.7b was used in this paper, while GPT-2 was used in (Miura et al., 2024b).

| Methods | Relevance | | | | Diversity | | Style |
|---|---|---|---|---|---|---|---|
| | BLEU-1 | BLEU-2 | ROUGE-1 | ROUGE-2 | Dist-1 | Dist-2 | StyCor |
| Ours-SCL (w/o style control loss) | 0.2105 | 0.1175 | 0.2954 | 0.1697 | 0.879 | 0.889 | 0.200 |
| Ours-II (w/o intimacy interpreter) | 0.2170 | 0.1257 | 0.3086 | 0.1826 | 0.907 | 0.917 | 0.247 |
| Ours$_{gold}$ | **0.2544** | **0.1463** | **0.3390** | **0.1999** | **0.925** | **0.930** | **0.250** |

Table 4: Results of Ablation Study

tomatic and human evaluations demonstrated that the proposed method outperformed the baseline in generating responses in a casual style when the user's level of intimacy was high and in a polite style when it was low.

The proposed dialog model was trained using a dialog corpus annotated with the speaker's level of intimacy. However, the availability of such a corpus is rather limited, while the construction of new corpora requires considerable costs. Therefore, it is essential to explore ways to enable LLMs to acquire the ability to control the style without relying on the intimacy-labeled corpus. Another important future work is to explore new style control frameworks that do not rely on pre-training the style language models and/or the style discrimination model.

## References

Noora Aapakallio. 2021. *Understanding Through Politeness – Translations of Japanese Honorific Speech to Finnish and English*. University of Eastern Finland.

Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. 2019. BERT: Pre-training of Deep Bidirectional Transformers for Language Understanding. In *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers)*, pages 4171–4186, Minneapolis, Minnesota. Association for Computational Linguistics.

Emily Dinan, Varvara Logacheva, Valentin Malykh, Alexander Miller, Kurt Shuster, Jack Urbanek, Douwe Kiela, Arthur Szlam, Iulian Serban, Ryan Lowe, et al. 2020. The second conversational intelligence challenge (convai2). In *The NeurIPS'18 Competition: From Machine Learning to Intelligent Conversations*, pages 187–208. Springer.

Joseph L. Fleiss and Cohen Jacob. 1973. The Equivalence of Weighted Kappa and the Intraclass Correlation Coefficient as Measures of Reliability. *Educational and Psychological Measurement*, 33(3):613–619.

Xiang Gao, Yizhe Zhang, Sungjin Lee, Michel Galley, Chris Brockett, Jianfeng Gao, and Bill Dolan. 2019.

Structuring Latent Spaces for Stylized Response Generation. In *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)*, pages 1814–1823, Hong Kong, China. Association for Computational Linguistics.

Ryuichiro Higashinaka, Kotaro Funakoshi, Michimasa Inaba, Yuiko Tsunomori, Tetsuro Takahashi, and Reina Akama. 2021. *Dialogue System Live Competition: Identifying Problems with Dialogue Systems Through Live Event*, pages 185–199. Springer Singapore.

Eduard Hovy. 1987. Generating natural language under pragmatic constraints. *Journal of Pragmatics*, 11(6):689–719.

Edward J Hu, Yelong Shen, Phillip Wallis, Zeyuan Allen-Zhu, Yuanzhi Li, Shean Wang, Lu Wang, Weizhu Chen, et al. 2022. Lora: Low-rank adaptation of large language models. *ICLR*, 1(2):3.

Yukiko Kageyama, Yuya Chiba, Takashi Nose, and Akinori Ito. 2018. Improving User Impression in Spoken Dialog System with Gradual Speech Form Control. In *Proceedings of the 19th Annual SIGdial Meeting on Discourse and Dialogue*, pages 235–240, Melbourne, Australia. Association for Computational Linguistics.

Chandra Khatri, Anu Venkatesh, Behnam Hedayatnia, Raefer Gabriel, Ashwin Ram, and Rohit Prasad. 2018. Alexa Prize — State of the Art in Conversational AI. *AI Magazine*, 39(3):40–55.

Kai Konen, Sophie Jentzsch, Diaoulé Diallo, Peer Schütt, Oliver Bensch, Roxanne El Baff, Dominik Opitz, and Tobias Hecking. 2024. Style Vectors for Steering Generative Large Language Models. In *Findings of the Association for Computational Linguistics: EACL 2024*, pages 782–802, St. Julian's, Malta. Association for Computational Linguistics.

Jinpeng Li, Zekai Zhang, Quan Tu, Xin Cheng, Dongyan Zhao, and Rui Yan. 2024. Stylechat: Learning recitation-augmented memory in llms for stylized dialogue generation. *arXiv preprint arXiv:2403.11439*.

Jiwei Li, Michel Galley, Chris Brockett, Jianfeng Gao, and Bill Dolan. 2016. A Diversity-Promoting Objective Function for Neural Conversation Models. In *Proceedings of the 2016 Conference of the North*

*American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pages 110–119, San Diego, California. Association for Computational Linguistics.

Chin-Yew Lin. 2004. ROUGE: A Package for Automatic Evaluation of Summaries. In *Text Summarization Branches Out*, pages 74–81, Barcelona, Spain. Association for Computational Linguistics.

Muxuan Liu and Ichiro Kobayashi. 2022. Construction and Validation of a Japanese Honorific Corpus Based on Systemic Functional Linguistics. In *Proceedings of the Workshop on Dataset Creation for Lower-Resourced Languages within the 13th Language Resources and Evaluation Conference*, pages 19–26, Marseille, France. European Language Resources Association.

Takuto Miura, Kiyoaki Shirai, Hideaki Kanai, and Natthawut Kertkeidkachorn. 2024a. Construction of a Japanese Dialog Corpus Annotated with Speakers' Intimacy. In *The 38th Pacific Asia Conference on Language, Information and Computation (PACLIC 38)*, page 10, Tokyo, Japan. Association for Computational Linguistics.

Takuto Miura, Kiyoaki Shirai, and Natthawut Kertkeidkachorn. 2024b. Intimacy-aware Style Control in Dialog Response Generation. In *The 9th Linguistic and Cognitive Approaches to Dialog Agents Workshop (LACATODA 2024)*, pages 5–16, Kyoto, Japan. CEUR-WS.

Tong Niu and Mohit Bansal. 2018. Polite Dialogue Generation Without Parallel Data. *Transactions of the Association for Computational Linguistics*, 6:373–389.

Kishore Papineni, Salim Roukos, Todd Ward, and Wei-Jing Zhu. 2002. BLEU: a method for automatic evaluation of machine translation. In *Proceedings of the 40th Annual Meeting on Association for Computational Linguistics*, ACL '02, page 311–318, USA. Association for Computational Linguistics.

Ashwin Ram, Rohit Prasad, Chandra Khatri, Anu Venkatesh, Raefer Gabriel, Qing Liu, Jeff Nunn, Behnam Hedayatnia, Ming Cheng, Ashish Nagar, et al. 2018. Conversational AI: The science behind the alexa prize. *arXiv preprint arXiv:1801.03604*.

Michael Silverstein. 2003. Indexical Order and the Dialectics of Social Life. *Language & Communication*, 23:193–229.

Mayumi Usami, editor. 2021. *BTSJ-Japanese Natural Conversation Corpus with Transcripts and Recordings (March 2021)*. National Institute for Japanese Language and Linguistics, Japan.

Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Ł ukasz Kaiser, and Illia Polosukhin. 2017. Attention is All you Need. In *Advances in Neural Information Processing Systems*, volume 30. Curran Associates, Inc.

Ronald Wardhaugh and Janet M Fuller. 2021. *An introduction to sociolinguistics*. John Wiley & Sons.

Ze Yang, Wei Wu, Can Xu, Xinnian Liang, Jiaqi Bai, Liran Wang, Wei Wang, and Zhoujun Li. 2020. StyleDGPT: Stylized Response Generation with Pretrained Language Models. In *Findings of the Association for Computational Linguistics: EMNLP 2020*, pages 1548–1559, Online. Association for Computational Linguistics.

Yinhe Zheng, Zikai Chen, Rongsheng Zhang, Shilei Huang, Xiaoxi Mao, and Minlie Huang. 2021. Stylized dialogue response generation using stylized unpaired texts. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 35, pages 14558–14567, Online. AAAI Press.

Qingfu Zhu, Weinan Zhang, Ting Liu, and William Yang Wang. 2021. Neural stylistic response generation with disentangled latent variables. In *Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing (Volume 1: Long Papers)*, pages 4391–4401, Bangkok, Thailand. Association for Computational Linguistics".

# Next Speaker Prediction for Multi-Speaker Dialogue with Large Language Models

**Lukas Hilgert**    **Jan Niehues**

Karlsruhe Institute of Technology, Germany

{lukas.hilgert, jan.niehues}@kit.edu

## Abstract

Currently Large Language Models (LLMs) are mostly used through a chatbot interface with the user manually deciding when the system should respond. In multi-speaker conversations (e.g., two humans and one robot) it is not clear who speaks when. We therefore investigate the ability of LLMs to predict the dialog structure. First, we frame the task as Next Speaker Prediction (NSP) and create a multi-domain test set. Secondly, we build dedicated systems for the NSP task using LLMs and finally performed automatic and human evaluation. Our final system matches the human performance when tested on unseen data and exceeds it on data of the same domain as the training data.

## 1 Introduction

In multi-speaker dialogues, it is important for the participants to know when to speak, as talking at the wrong time may be irritating for the other speakers and may even hinder the speakers to reach their goals. It is crucial for dialogue systems to handle this task well as speaking too often may be annoying to the user while speaking rarely may seem unresponsive to the user and opposes the system's purpose.

Large Language Models (LLMs) are the core of modern dialogue systems. Currently they are mostly used through a chatbot interface where they only respond after the user sends a chat message. Here, there is no need for dedicated dialog structure modeling as the user always decides when the model should respond. For spoken dialogue with two speakers, the modeling is not as trivial as it is not clear when one speaker ends their turn. For multi-speaker scenarios it is significantly more challenging when the LLM should respond as the users could be chatting with each other directly during the course of dialogue.

P11: i used to live downtown san jose and every once in a while i just get with garlic and i don't know if it's from gilroy probably not nut i like to think it was so [laugh]
P09: yeah
P09: wow
P12: what are actually some nice places to go around here cause i've moved here recently so [unintelligible]
P09: napa napa is nice
P10: oh
P10: napa is nice [unintelligible]
P12: oh yeah actually i went there last week and they had uhhh i think sonoma had a hot air balloon festival there [...] but it's pretty nice seeing them at sunrise so yeah it was really beautiful yeah

Annotators' votes:
P09: 7, P11: 2, P10: 2

Zero-shot LLM: P12
Fine-tuned LLM: P09

Next utterance:
P09: people like to go wine tasting

Figure 1: Example of a part of a dialogue from DiPCo (Segbroeck et al., 2020). We show the previous utterances, which next speaker our human annotators predicted, what the LLMs in different setting predicted, and what the actual next speaker and utterance are.

We model this ability as the Next Speaker Prediction (NSP) task like Wei et al. (2023). We think it is a suitable proxy task as good performance on predicting the next speaker should indicate the quality of the system's ability to decide the correct time to actively contribute to the conversation.

We want to investigate multi-speaker dialogue from multiple domains to test generalization, estimate the performance of LLMs, and find out how well they have to perform. Therefore, our research on the NSP task covers the following aspects:

- We create a multi-domain benchmark for the NSP task utilizing multiple existing dialogue datasets.
- We run a user study with eleven annotators to gather a human baseline. This evaluation gives insights on the ambiguity of the task.

- We analyze the ability of various size LLMs to perform the NSP task and build dedicated models that reach or exceed our estimate of human performance.

## 2  Next Speaker Prediction Benchmark

To evaluate how well our approaches perform on the NSP task, we compile a benchmark consisting of multiple datasets. Using datasets from multiple domains enable us to estimate the generalization of the evaluated systems. Additionally, we collect a human baseline on subsamples of the datasets to get an estimate of the human performance on NSP. While the dialogue structure from the datasets offers a ground truth for the next speaker, we want to find out if human annotators would consider other options as equally possible. Also, we obtain an overview how ambiguous the task is for human annotators.

### 2.1  Datasets

For the NSP task, we need datasets of dialogues where the speaker is denoted for every utterance. As the following utterance then always determines which speaker will be the next after the current one, we can easily model the NSP task. For dialogues with only two speakers, the NSP task is fairly trivial. Therefore, we only investigated multi-speaker dialogue datasets.

We use three dialogue datasets for our benchmark (Table 1) to cover multiple domains. We chose these three datasets to cover multiple domains. Also, there is an existing baseline for the NSP task for MultiLIGHT (ML, Wei et al. (2023)). The other two datasets both include the type of noise that a dialogue system would also encounter in a real-life setting. Additionally, the conversation domains are realistic settings than ML's (Table 1). We use two of these similar datasets as this allows us to compare how well our approaches generalize an unseen domain and different noise as DiPCo includes no training data. The datasets differ in the numbers of participating speakers in one conversation, the domain of the conversation (topic, setting), and the amount of noise in the sense of very short utterances that introduce no or almost no substance to the conversation.

ML is a text-only dataset created specifically for dialogue research. The authors also performed experiments on the NSP task with at time of publication current Transformer-based language models

| Dataset | AMI | DiPCo | ML |
|---|---|---|---|
| # Speakers | 4 | 4 | 3 |
| Domain | meeting | dinner party | fantasy role-play |
| Noisy | yes | yes | no |
| # Utterance | 12627 | 3400 | 9164 |
| # Dialogues | 16 | 5 | 323 |
| Avg. utts. | 789.19 | 680.00 | 28.37 |

Table 1: Properties of the investigated datasets (specific numbers from the test splits). We list the number of speakers per dialogue, the topics of the conversations, if they contain some form of noise (short / interrupting utterances), and the number of utterances in total, the number of dialogues, and the average number of consecutive utterances per dialogue.

| Dataset | AMI | DiPCo |
|---|---|---|
| Speaker 0 | 32.18 | 23.93 |
| Speaker 1 | 26.88 | 25.75 |
| Speaker 2 | 23.36 | 28.04 |
| Speaker 3 | 18.75 | 22.28 |

Table 2: Contributed utterances (in percentage) from each speaker across all dialogues. For AMI, the speaker that speaks earlier in the dialogue, seems to have more dialogue utterances while there seems to be now such accumulation for DiPCo.

that they fine-tuned on this task. The AMI meeting corpus (Carletta et al., 2005) and the Dinner Party Corpus (DiPCo) (Segbroeck et al., 2020) are primarily audio (and video for AMI) datasets from recorded conversations.

The type of conversations in AMI are meetings and in DiPCo dinner party talk. Both contain noise like "Umm", "Hmm", and "Yeah" that introduce no or almost no substance to the conversation in some cases. While these appear to happen at random times, these kinds of utterances are also present in a setting where an LLM gets its input via an Automatic Speech Recognition system. Also, for utterances like "Yeah" it is hard to determine if "Yeah" is just noise or an import acknowledgment of a previous utterance. So, we only filter out obvious irrelevant utterances for the DiPCo dataset like "[Noise]" to reduce the noisiness while keeping potentially important utterances.

**Datasets Statistics**  In a first step, we investigated the dataset statistics in order to identify the various challenges of the datasets.

For example, we analyzed the percentage of con-

| Dataset | AMI | DiPCo | ML |
|---|---|---|---|
| 4 | 91.51 | 89.66 | 29.32 |
| 8 | 66.79 | 49.33 | 14.04 |
| 16 | 40.53 | 20.42 | 13.48 |
| 32 | 20.71 | 6.04 | 13.48 |
| 64 | 8.60 | 0.84 | 13.48 |

Table 3: Percentage of contexts where at least one speaker is missing depending on the number of recent utterances included in the prompt.

tributed utterances per speaker within each dialogue to see if one specific speaker speaks significantly more often which could lead to a bias to predict that speaker more often as the next one. We number the speakers ascending by their order of appearance. For AMI, the speakers that appear earlier in the conversation seem to speak more often. After qualitative analysis, we concluded that this is the case because in AMI the person opening the meeting is also the organizer of the meeting itself. We saw no such clear trend for DiPCo.

We want to only include the recent dialogue utterances in our benchmark as the dialogues in the datasets are up to several hundred utterances long (Table 1) which could overwhelm both human annotators and NLP systems. We therefore examined the number of times where at least one speaker is missing from our dialogue excerpt to find out in how many cases the context is missing information about some speakers. We start with four included recent utterances and iteratively double the amount up to 64. For ML, the number does not continue to decrease after 16 included utterances (Table 3). This is a result of the fact that in the beginning of the dialogue, not all speakers have spoken yet. As the dialogues in ML are short, this situation is quite common. For the other two, including quadratically more recent utterances linearly reduces the number of excerpts with missing speakers. This shows that very often in a small enough context window only a subset of the speakers interact with each other.

## 2.2 Human Baselines

In a first step, we analyze the difficulty of the task through a human evaluation. While the dialogues from the datasets were generated by humans, like many other Natural Language Processing (NLP) tasks, the NSP task is also ambiguous. We therefore collect human data on the NSP task for samples of consecutive utterances of the test splits of all three

datasets. Our sample size is 63 dialogue utterances for AMI (0.50% of the full test set), 55 for DiPCo (2.00%), and 91 for ML (0.96%). As the dialogues in ML are fairly short, our sample includes three full dialogues. These sample sizes should in our opinion capture the natures of the datasets while also keeping the annotation work at a reasonable level. The user study involved eleven participants for each dataset. We average each's accuracy to get the human baseline (Table 6).

We included the last 32 utterances and did not rename the speakers in the prompts. We chose 32 as this number is higher than the number of utterances in full dialogues for the ML dataset and is not overwhelmingly large for human annotators. For the names, we assumed that the annotators should be able to distinguish the names more easily with the original ones from the dataset.

| Dataset | Fleiss' kappa |
|---|---|
| AMI | 0.17 |
| DiPCo | 0.14 |
| ML 1 | 0.49 |
| ML 2 | 0.43 |
| ML 3 | 0.32 |

Table 4: Fleiss' kappa for multi-rater agreement on the samples used for the gathering the human baseline.

We provide the Fleiss' kappa multi-rater agreement measure (Fleiss, 1971) for each dataset sample (Table 4). For ML, we show the score for each of three dialogues that are included in our sample. The scores low showing the ambiguity of the task. The difference between ML and the other two datasets are in our opinion a result of it having fewer speakers per dialogue and having less noisy utterances. Manual inspection and anecdotal evidence from the annotators showed that the annotators agreed or were sure in their prediction respectively for some turns (most annotators picked one speaker) but disagreed or were unsure in their prediction respectively in other cases (annotators picked different speakers, no clear "favorite").

## 3 Next Speaker Prediction with LLMs

We want to use state-of-the-art technology to build a next speaker predictor. This leads to LLMs as they excel on other NLP tasks. Additionally, their task during the pre-training phase is predicting the next token which corresponds to predicting the next speaker when the prompt is a dialogue transcript

with annotated speakers. This implies that the NSP task is "natural" for LLMs given their training.

While the authors of ML perform similar experiments, they were with the smaller encoder-decoder language models R2C2 (Shuster et al., 2022), T5 (Raffel et al., 2020), and BART (Lewis et al., 2020) which are smaller than today's models and were trained on less data and did not receive the extensive post-training of current LLMs. Furthermore, the authors of ML had to fine-tune these models to perform this task while current LLMs can be used with zero-shot prompts.

To model the NSP task as an LLM task, we prompt the LLMs to predict the next speaker by utilizing the information we provide (Appendix A): An instruction for the task and the most recent utterances of the current dialogue as context. Each utterance starts with the corresponding speaker. Each dataset already contains identifiers for the speakers. For ML, each speaker has a descriptive name like "jester". The other two datasets use string identifiers like "P12" or "MTD011UID".

For every dialogue turn we include the last eight utterances as context for the LLMs and rename the speakers to the same generic identifier across all datasets to increase the similarity of the task across the datasets. We replace them with renamings where each speaker has a pseudonym in the format of "speaker <number>". ML also includes descriptions of the character of each speaker and the location of the dialogue. We did not include this information to keep the task comparable.

Although LLMs are able to perform on zeroshot, often specialized models perform better. We therefore train LLMs supervised on the NSP task on multi-domain data by mixing the training splits from the AMI and the ML dataset. We use a balanced mixture (similar number of training data points) to ensure generalization across domains. We train with pairs of the prompt used in the zeroshot setting and the expected speaker from the datasets, so that the model learns how to map the recent dialogue turns to the next speaker.

## 4 Experiments

We evaluate the LLMs on the test splits of the datasets, compare them to random and human baselines, and perform ablation studies on our modeling decisions.

We chose next speaker accuracy as our main evaluation metrics as this is the most straightforward

metric with the given data. As the distributions of utterances per speaker are fairly balanced (Table 2), we did not employ metrics like $F_1$. While accuracy is a "hard" metric and does not account for ambiguity, we assume that the fairly large dataset size and direct comparison against baselines still gives a good estimate how well the LLMs (and especially our fine-tuned one) do for NSP. Nevertheless, we analyze the agreement of the LLMs with the annotators (section 4.3).

### 4.1 Setup

We perform all our experiments with models from the Llama 3 family (Dubey et al., 2024). We use the 3B (3.2. 3B) and 8B (3.1 8B) parameter version for zero-shot and fine-tuning experiments while we use the bigger version (3.3 70B) only in a zero-shot setting as fine-tuning this model requires significantly more compute and the smaller models responded already very well to fine-tuning.

The fine-tuning data mixture consists of all the available training data from the AMI meeting corpus and 33% from ML. We use only 33% to balance the number of data points from each dataset. DiPCo has no train split. We conduct ablation studies on all mentioned modeling decisions including the preprocessing (subsection 4.4). We made these decisions that impacted our main results during development on the basis of the validation sets which all utilized datasets provide.

### 4.2 Random Baselines

To compare our results to another baseline, we present three random baselines. Each is designed to model two very distinct types of dialogue flow and a combination of both. These baselines are: One where the speaker is picked randomly but always switches after each dialogue utterance (denoted as **always**). Then, we assume that the speaker never switches, so we predict the last speaker to also be the next speaker (denoted as **never**). At last, we model a combination of both where we pick the speaker completely randomly without excluding the last speaker (denoted as **usually**). We run each method five times and average the results.

### 4.3 Main Experiments

We differentiate between the results on the full test sets and the samples for the human baseline.

**Results on the full Tests Sets** The accuracy scores (Table 5) for the random baselines illustrate

| Dataset | AMI | DiPCo | ML |
|---|---|---|---|
| Random Baselines | | | |
| Always | 22.10 | 26.62 | 45.91 |
| Usually | 25.17 | 25.21 | 33.36 |
| Never | 33.41 | 19.32 | 8.91 |
| Zero-shot | | | |
| Llama 3.2 3B | 25.28 | 25.66 | 28.10 |
| Llama 3.1 8B | 34.88 | 30.94 | 40.41 |
| Llama 3.3 70B | 35.81 | 32.98 | 52.06 |
| Fine-tuned | | | |
| Llama 3.2 3B | 45.91 | 36.91 | 59.40 |
| Llama 3.1 8B | 47.85 | 38.48 | 59.85 |

Table 5: NSP accuracy on the **full test splits**. We compare the accuracy of the random baselines and the Llama 3 models in a zero-shot and fine-tuned setting. Fine-tuning improves performance beyond the 70B model's performance. Even the dataset we did not train on (DiPCo) benefits from fine-tuning on the NSP task.

| Dataset | AMI | DiPCo | ML |
|---|---|---|---|
| Human | 30.88 | 33.22 | 48.65 |
| Random Baselines | | | |
| Always | 20.63 | 27.64 | 45.49 |
| Usually | 17.78 | 26.91 | 35.16 |
| Never | 32.06 | 14.55 | 11.87 |
| Zero-shot | | | |
| Llama 3.2 3B | 15.87 | 21.82 | 25.27 |
| Llama 3.1 8B | 34.92 | 23.64 | 32.97 |
| Llama 3.3 70B | 30.16 | 34.55 | 51.65 |
| Fine-tuned | | | |
| Llama 3.2 3B | 47.62 | 40.00 | 61.54 |
| Llama 3.1 8B | 58.73 | 34.55 | 59.34 |

Table 6: NSP accuracy on the samples of the test splits for the **human baseline**. We compare the accuracy of the human annotators, random baselines, and the Llama 3 models in a zero-shot and fine-tuned setting. Fine-tuning beats human accuracy on the datasets with training data but also on DiPCo.

what we already saw during qualitative analysis of the datasets: In the AMI meeting corpus, the speakers often deliver multiple utterances after another while in the ML dataset the speaker almost always switches. Llama 3.1 8B performs a bit or clearly better than the random baselines on AMI and DiPCo, which highlights the importance of a multi-domain benchmark. On ML however, simply randomly picking one of the other two speaker as the next performs better. The smallest model we tested (3.2 3B) only manages to predict next speaker as well as completely randomly picking one. The bigger 70B model outperforms the random baselines clearly on DiPCo and ML. We see a clear trend that scaling the model size increases the ability to predict the next speaker.

When fine-tuning 3.1 8B on the task, it significantly outperforms itself in the zero-shot setting, the random baselines, and the bigger version. The performance even improves beyond the 70B model on the DiPCo dataset, which has no training split meaning that this dataset is out-of-domain for the fine-tuned models, and we see generalization for different domains. The case for 3.2 3B is similar but with slightly lower scores than 3.1 8B.

**Results on the Samples of Tests Sets for Human Baselines** On DiPCo and ML, our collected human baseline outperforms the random baselines albeit not all of them by a big margin (Table 6). For AMI, it is even slightly below the best technique ("never") that assumes the last speaker will always be the next speaker. In the zero-shot setting, the smallest Llama model shows the same pattern as on the full test sets. The medium LLM however achieves a higher accuracy on AMI as the human baseline, while struggling to reach the random baseline on the other two datasets which may be specific to these samples. The 70B version roughly matches the human performance on all datasets. The scaling trends we observed on the full test sets is also present on the samples except for AMI, where the 70B model underperforms the 8B model.

Fine-tuning the two smaller models shows similar effects as we saw on the full test split: The NSP accuracy is increased greatly compared to the zero-shot setting and even sightly outperforms the 70B model on the datasets where training data exists. For DiPCo, the performance of Llama 3.1 8B is the same as the one of 3.3 70B. The fine-tuned 3B model manages to outperform both the 8B and 70B model on DiPCo. As it showed reduced performance compared to the 8B model on the full test sets and as this sample set is small, we assume that these differences between the models are partly noise while still showing the effectiveness of our fine-tuning in general for the NSP task.

**Agreement of Annotators and LLMs** As mentioned before, this task is a highly ambiguous task. However, there are also situations where only a small set of possible next speakers are correct. We

wanted to investigate this and therefore use the human annotations as additional references.

We analyze the agreement of the LLMs with the human annotators. To do this, we remove one annotator at a time from the pool of annotators. We then compare their agreement with the rest of the annotators and with the LLMs by measuring the accuracy of their predictions. We then average the results for all annotators.

Additionally, we show how many of the predictions can be counted as correct with these conditions which decreases with the number of required agreeing annotators increasing (row "Correct answers", column "all").

In this setup, we counted a prediction as correct if at least $n$ annotators propose this prediction. This allows for situations where then no answer is correct and therefore it does not matter what the model predicts and for situation where multiple solutions are correct. Additionally, we show how many of the predictions can be counted as correct with these conditions which decreases with the number of required agreeing annotators increasing (row "Correct answers", column "all").

Also, we show how many choices a predictor has with the given threshold as for example only three possible next speakers can be counted as correct if the number of annotators is ten and the threshold for the number of agreeing annotator is three. Therefore, the number of possible correct answers also decreases with a higher threshold. The reported numbers for the annotators and the models display the percentage of correct predictions (given a threshold) out of the possible correct answers. We then also list the distribution of choices within this set – how many predictions are possibly correct. Per bin of possibly correct prediction, we also report the accuracy of each predictor.

For the AMI dataset, we see mixed results: From a threshold of three and more, the larger model has lower agreement than the 8B model. The fine-tuned model shows a similar regression for a threshold of three and five. For seven agreeing annotators, the fine-tuned model has a slightly higher agreement, yet the 70B model is lower than Llama 3.1 8B in zero-shot. We think that these results come from the fact that fine-tuning on the AMI training data pushed the 8B LLM towards the distribution by the dataset increasing the NSP accuracy, which disagrees with our human annotators. That the 70B model also has a lower agreement could be a sign of its training data containing part of AMI and it

memorizing it better than the 8B model.

For DiPCo, we see that the 8B model in the fine-tuned setting has a clearly higher (threshold of one and three) or slightly higher (threshold of five) agreement than in the zero-shot setting (Table 8). Here, we also see that the 70B version has higher agreement than the 8B model in zero-shot. This matches our observations from the accuracy scores before that increased model sizes correlates with an improved NSP ability. Fine-tuning Llama 8B therefore improves for most tested thresholds the agreement with the human annotators on DiPCo and moves it closer to that of the 70B model. As we did not fine-tune the 8B model on data from DiPCo, we think that these results together with the increase in NSP accuracy show that training on the NSP task with dialogue datasets does generalize to better NSP performance – matching or exceeding human performance in NSP accuracy.

## 4.4   Ablation Studies

We also examine our modeling decisions when fine-tuning Llama 3.1 8B.

| Dataset | AMI | DiPCo | ML |
|---|---|---|---|
| Speaker Renaming | | | |
| Original | 42.04 | 39.35 | 54.58 |
| **Renamed** | 47.85 | 38.48 | 59.85 |
| Context Length | | | |
| 4 | 46.32 | 34.66 | 59.47 |
| **8** | 47.85 | 38.48 | 59.85 |
| 16 | 47.92 | 39.46 | 60.13 |
| 32 | 47.58 | 37.86 | 60.21 |
| 64 | 46.72 | 36.29 | 59.73 |
| Training Data Mixture | | | |
| Zero-shot | 34.88 | 30.94 | 40.41 |
| AMI | 47.84 | 37.35 | 42.75 |
| ML | 24.08 | 28.25 | 60.07 |
| **AMI + 33% ML** | 47.85 | 38.48 | 59.85 |

Table 9: Comparison of the accuracy results from the ablation studies. Renaming the speakers to a dataset-across scheme increases performance in general. Including more previous utterances in the prompt only helps until 16 utterances. Training only on one of the two available datasets is worse than using both.

**Speaker Renaming**   We compare the unmodified versions of the datasets with our renamed versions. Renaming improves performance on all datasets except for DiPCo (Table 9). This is probably the case as the speaker names in DiPCo (e.g., "P12") are already fairly generic but distinct. This also

| # choices | all | 1 | 2 | 3 | 4 |
|---|---|---|---|---|---|
| At least one out of ten agreeing annotator | | | | | |
| Correct answers | 100.00 | 2.31 | 21.07 | 43.00 | 33.62 |
| Annotators | 92.06 | 68.75 | 86.99 | 89.60 | 100.00 |
| Zero-shot 8B | 88.46 | 75.00 | 63.70 | 92.28 | 100.00 |
| Fine-tuned 8B | 88.74 | 81.25 | 87.67 | 80.87 | 100.00 |
| Zero-shot 70B | 89.32 | 75.00 | 69.18 | 91.61 | 100.00 |
| At least three of ten agreeing annotator | | | | | |
| Correct answers | 100.00 | 42.14 | 54.98 | 2.89 | 0.00 |
| Annotators | 68.40 | 56.51 | 76.12 | 95.00 | 0.00 |
| Zero-shot 8B | 66.67 | 43.49 | 82.68 | 100.00 | 0.00 |
| Fine-tuned 8B | 54.83 | 44.86 | 61.42 | 75.00 | 0.00 |
| Zero-shot 70B | 61.18 | 34.93 | 79.27 | 100.00 | 0.00 |
| At least five of ten agreeing annotator | | | | | |
| Correct answers | 82.40 | 98.42 | 1.58 | 0.00 | 0.00 |
| Annotators | 56.39 | 56.23 | 66.67 | 0.00 | 0.00 |
| Zero-shot 8B | 45.01 | 44.13 | 100.00 | 0.00 | 0.00 |
| Fine-tuned 8B | 39.93 | 39.15 | 88.89 | 0.00 | 0.00 |
| Zero-shot 70B | 39.23 | 39.32 | 33.33 | 0.00 | 0.00 |
| At least seven of ten agreeing annotator | | | | | |
| Correct answers | 29.29 | 100.00 | 0.00 | 0.00 | 0.00 |
| Annotators | 59.61 | 59.61 | 0.00 | 0.00 | 0.00 |
| Zero-shot 8B | 44.33 | 44.33 | 0.00 | 0.00 | 0.00 |
| Fine-tuned 8B | 45.81 | 45.81 | 0.00 | 0.00 | 0.00 |
| Zero-shot 70B | 42.36 | 42.36 | 0.00 | 0.00 | 0.00 |

Table 7: Agreement between annotators and LLMs (**AMI**): We show the NSP accuracy for each annotator (results averaged) and the LLMs when the other annotators serve as the ground truth. We show different thresholds for agreeing annotators that an answer counts as correct. We also display the accuracy grouped by the number of choices a predictor has (if too many annotators have to agree, the number of possible correct answers shrink).

means that not renaming the speakers for the user study should not skew our comparison.

**Context Length**   We also compare how the number of included most recent dialogue utterances influences the accuracy of the predictions: We vary the number of included utterances in the prompt as context for the models in steps of the power of two from four to 64. There seems to be a limit on how much context in the form of previous dialogue utterances helps the model in its decision even with the number of not included speakers decreasing (Table 3). We picked eight recent utterances for our experiments as it showed the best performance on the validation sets, and it enables faster inference than for 16 utterances. As the accuracies differ only sightly across the context lengths we tried, it seems that the model mostly relies on the last few utterances for its decision while also being able to focus on them even if the included dialogue is longer.

**Training Data Mixture**   As we have two datasets from our benchmark with training data, we want to find out how the specific selection of training data impacts the generalization ability of the fine-tuned models. Only training on the AMI data already shows large improvements for the two similar datasets (AMi and DiPCo) but only small improvements for ML. Only training on this dataset however reduces the performance on the other two datasets. A weighted combination of both datasets (roughly equal amount of datapoints from both) resulted in performance similar like training on the "corresponding" dataset. We even saw slight transfer learning for DiPCo.

## 5   Related Work

Previous research on dialogue turns is different from our approach as we assume both the setting of a multi-speaker dialogue either in text form or as a transcript and a text-only LLM as the predictor.

| # choices | all | 1 | 2 | 3 | 4 |
|---|---|---|---|---|---|
| At least one of ten agreeing annotator | | | | | |
| Correct answers | 100.00 | 0.00 | 5.95 | 50.58 | 43.47 |
| Annotators | 92.07 | 0.00 | 61.11 | 88.89 | 100.00 |
| Zero-shot 8B | 92.56 | 0.00 | 88.89 | 86.60 | 100.00 |
| Fine-tuned 8B | 98.18 | 0.00 | 94.44 | 97.06 | 100.00 |
| Zero-shot 70B | 96.53 | 0.00 | 61.11 | 97.71 | 100.00 |
| At least three of ten agreeing annotator | | | | | |
| Correct answers | 100.00 | 44.30 | 51.74 | 3.97 | 0.00 |
| Annotators | 61.98 | 54.85 | 67.73 | 66.67 | 0.00 |
| Zero-shot 8B | 55.21 | 37.69 | 67.41 | 91.67 | 0.00 |
| Fine-tuned 8B | 61.98 | 50.00 | 70.93 | 79.17 | 0.00 |
| Zero-shot 70B | 61.98 | 48.13 | 72.52 | 79.17 | 0.00 |
| At least five of ten agreeing annotator | | | | | |
| Correct answers | 70.74 | 100.00 | 0.00 | 0.00 | 0.00 |
| Annotators | 55.61 | 55.61 | 0.00 | 0.00 | 0.00 |
| Zero-shot 8B | 46.26 | 46.26 | 0.00 | 0.00 | 0.00 |
| Fine-tuned 8B | 46.73 | 46.73 | 0.00 | 0.00 | 0.00 |
| Zero-shot 70B | 53.04 | 53.04 | 0.00 | 0.00 | 0.00 |
| At least seven of ten agreeing annotator | | | | | |
| Correct answers | 25.12 | 100.00 | 0.00 | 0.00 | 0.00 |
| Annotators | 43.42 | 43.42 | 0.00 | 0.00 | 0.00 |
| Zero-shot 8B | 45.39 | 45.39 | 0.00 | 0.00 | 0.00 |
| Fine-tuned 8B | 42.76 | 42.76 | 0.00 | 0.00 | 0.00 |
| Zero-shot 70B | 50.66 | 50.66 | 0.00 | 0.00 | 0.00 |

Table 8: Agreement between annotators and LLMs (**DiPCo**): We show the NSP accuracy for each annotator (results averaged) and the LLMs when the other annotators serve as the ground truth. We show different thresholds for agreeing annotators that an answer counts as correct. We also display the accuracy grouped by the number of choices a predictor has (if too many annotators have to agree, the number of possible correct answers shrink).

**Transition Relevance Places** Methods for turn-taking use LLMs to predict transition-relevant places within a stream of words. Transition-relevant places are points in a dialogue where a turn-shift can happen. Ekstedt and Skantze (2020) fine-tuned GPT-2 to predict these spots in written and spoken dialogues. Later work (Umair et al., 2024) investigated if more recent LLMs (e.g., Llama 3.1 8B) can do the same.

**Audio / Visual Cues** Multimodal approaches for NSP use visual cues like gaze and hand gestures (Ishii et al., 2016; Malik et al., 2020). This research incorporates gaze transition patterns and eye contact timing structure (Ishii et al., 2016) or head movement (Ishii et al., 2015) to predict the next speaker using support vector machines. Malik et al. (2020) utilized focus of attention among others to train classic machine learning classifiers for NSP. Other systems rely on voice activity projection for turn-taking prediction (Inoue et al., 2024a,b) which

predicts future voice activity based on the current audio signal.

## 6 Conclusion

Our research goal was to investigate the ability of LLMs to predict the next speaker in a multi-speaker dialogue setting. We also compared their performance with humans and fine-tuned LLMs to improve them on NSP. The experiments on our compiled benchmark show that LLMs like Llama 3.3 70B can match the human performance on the NSP task in accuracy and it also shows very high agreement with human predictors. Smaller LLMs can achieve this performance or even exceed it by fine-tuning on dialogue datasets when the dialogue flow (e.g., with some short noisy utterances) is similar. We think that these results imply an ability of LLMs to "know" when to talk at transition-relevant places in a multi-speaker dialogue – either through large model size or fine-tuning on dialogues. Fu-

ture work will investigate how multimodal LLMs handle the NSP task as this work did not investigate the impact of additional auditory and visual information about the dialogue.

## Limitations

Our investigation is limited to text-only dialogues and does not cover the use of audio or visual cues. We do not predict the next speaker on a per-token or per-word basis but rather after a full utterance. This assume that the system only receives full utterances as input which is the case if the dialogue participants interact via text or through an audio transcript.

## Acknowledgments

## References

Jean Carletta, Simone Ashby, Sebastien Bourban, Mike Flynn, Maël Guillemot, Thomas Hain, Jaroslav Kadlec, Vasilis Karaiskos, Wessel Kraaij, Melissa Kronenthal, Guillaume Lathoud, Mike Lincoln, Agnes Lisowska, Iain McCowan, Wilfried M. Post, Dennis Reidsma, and Pierre Wellner. 2005. The AMI meeting corpus: A pre-announcement. In *Machine Learning for Multimodal Interaction, Second International Workshop, MLMI 2005, Edinburgh, UK, July 11-13, 2005, Revised Selected Papers*, volume 3869 of *Lecture Notes in Computer Science*, pages 28–39. Springer.

Abhimanyu Dubey, Abhinav Jauhri, Abhinav Pandey, Abhishek Kadian, Ahmad Al-Dahle, Aiesha Letman, Akhil Mathur, Alan Schelten, Amy Yang, Angela Fan, Anirudh Goyal, Anthony Hartshorn, Aobo Yang, Archi Mitra, Archie Sravankumar, Artem Korenev, Arthur Hinsvark, Arun Rao, Aston Zhang, Aurélien Rodriguez, Austen Gregerson, Ava Spataru, Baptiste Rozière, Bethany Biron, Binh Tang, Bobbie Chern, Charlotte Caucheteux, Chaya Nayak, Chloe Bi, Chris Marra, Chris McConnell, Christian Keller, Christophe Touret, Chunyang Wu, Corinne Wong, Cristian Canton Ferrer, Cyrus Nikolaidis, Damien Allonsius, Daniel Song, Danielle Pintz, Danny Livshits, David Esiobu, Dhruv Choudhary, Dhruv Mahajan, Diego Garcia-Olano, Diego Perino, Dieuwke Hupkes, Egor Lakomkin, Ehab AlBadawy, Elina Lobanova, Emily Dinan, Eric Michael Smith, Filip Radenovic, Frank Zhang, Gabriel Synnaeve, Gabrielle Lee, Georgia Lewis Anderson, Graeme Nail, Grégoire Mialon,

Guan Pang, Guillem Cucurell, Hailey Nguyen, Hannah Korevaar, Hu Xu, Hugo Touvron, Iliyan Zarov, Imanol Arrieta Ibarra, Isabel M. Kloumann, Ishan Misra, Ivan Evtimov, Jade Copet, Jaewon Lee, Jan Geffert, Jana Vranes, Jason Park, Jay Mahadeokar, Jeet Shah, Jelmer van der Linde, Jennifer Billock, Jenny Hong, Jenya Lee, Jeremy Fu, Jianfeng Chi, Jianyu Huang, Jiawen Liu, Jie Wang, Jiecao Yu, Joanna Bitton, Joe Spisak, Jongsoo Park, Joseph Rocca, Joshua Johnstun, Joshua Saxe, Junteng Jia, Kalyan Vasuden Alwala, Kartikeya Upasani, Kate Plawiak, Ke Li, Kenneth Heafield, Kevin Stone, and et al. 2024. The Llama 3 herd of models. *CoRR*, abs/2407.21783.

Erik Ekstedt and Gabriel Skantze. 2020. TurnGPT: a transformer-based language model for predicting turn-taking in spoken dialog. In *Findings of the Association for Computational Linguistics: EMNLP 2020, Online Event, 16-20 November 2020*, volume EMNLP 2020 of *Findings of ACL*, pages 2981–2990. Association for Computational Linguistics.

Joseph L Fleiss. 1971. Measuring nominal scale agreement among many raters. *Psychological bulletin*, 76(5):378.

Edward J. Hu, Yelong Shen, Phillip Wallis, Zeyuan Allen-Zhu, Yuanzhi Li, Shean Wang, Lu Wang, and Weizhu Chen. 2022. LoRA: Low-rank adaptation of large language models. In *The Tenth International Conference on Learning Representations, ICLR 2022, Virtual Event, April 25-29, 2022*. OpenReview.net.

Koji Inoue, Bing'er Jiang, Erik Ekstedt, Tatsuya Kawahara, and Gabriel Skantze. 2024a. Multilingual turn-taking prediction using voice activity projection. In *Proceedings of the 2024 Joint International Conference on Computational Linguistics, Language Resources and Evaluation, LREC/COLING 2024, 20-25 May, 2024, Torino, Italy*, pages 11873–11883. ELRA and ICCL.

Koji Inoue, Bing'er Jiang, Erik Ekstedt, Tatsuya Kawahara, and Gabriel Skantze. 2024b. Real-time and continuous turn-taking prediction using voice activity projection. *CoRR*, abs/2401.04868.

Ryo Ishii, Shiro Kumano, and Kazuhiro Otsuka. 2015. Predicting next speaker based on head movement in multi-party meetings. In *2015 IEEE International Conference on Acoustics, Speech and Signal Processing, ICASSP 2015, South Brisbane, Queensland, Australia, April 19-24, 2015*, pages 2319–2323. IEEE.

Ryo Ishii, Kazuhiro Otsuka, Shiro Kumano, and Junji Yamato. 2016. Prediction of who will be the next speaker and when using gaze behavior in multiparty meetings. *ACM Trans. Interact. Intell. Syst.*, 6(1):4:1–4:31.

Mike Lewis, Yinhan Liu, Naman Goyal, Marjan Ghazvininejad, Abdelrahman Mohamed, Omer Levy, Veselin Stoyanov, and Luke Zettlemoyer. 2020. BART: denoising sequence-to-sequence pre-training

for natural language generation, translation, and comprehension. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics, ACL 2020, Online, July 5-10, 2020*, pages 7871–7880. Association for Computational Linguistics.

Usman Malik, Julien Saunier, Kotaro Funakoshi, and Alexandre Pauchet. 2020. Who speaks next? turn change and next speaker prediction in multimodal multiparty interaction. In *32nd IEEE International Conference on Tools with Artificial Intelligence, IC-TAI 2020, Baltimore, MD, USA, November 9-11, 2020*, pages 349–354. IEEE.

Colin Raffel, Noam Shazeer, Adam Roberts, Katherine Lee, Sharan Narang, Michael Matena, Yanqi Zhou, Wei Li, and Peter J. Liu. 2020. Exploring the limits of transfer learning with a unified text-to-text transformer. *J. Mach. Learn. Res.*, 21:140:1–140:67.

Maarten Van Segbroeck, Ahmed Zaid, Ksenia Kutsenko, Cirenia Huerta, Tinh Nguyen, Xuewen Luo, Björn Hoffmeister, Jan Trmal, Maurizio Omologo, and Roland Maas. 2020. DiPCo - dinner party corpus. In *21st Annual Conference of the International Speech Communication Association, Interspeech 2020, Virtual Event, Shanghai, China, October 25-29, 2020*, pages 434–436. ISCA.

Kurt Shuster, Mojtaba Komeili, Leonard Adolphs, Stephen Roller, Arthur Szlam, and Jason Weston. 2022. Language models that seek for knowledge: Modular search & generation for dialogue and prompt completion. In *Findings of the Association for Computational Linguistics: EMNLP 2022, Abu Dhabi, United Arab Emirates, December 7-11, 2022*, pages 373–393. Association for Computational Linguistics.

Muhammad Umair, Vasanth Sarathy, and Jan Peter de Ruiter. 2024. Large language models know what to say but not when to speak. In *Findings of the Association for Computational Linguistics: EMNLP 2024, Miami, Florida, USA, November 12-16, 2024*, pages 15503–15514. Association for Computational Linguistics.

Jimmy Wei, Kurt Shuster, Arthur Szlam, Jason Weston, Jack Urbanek, and Mojtaba Komeili. 2023. Multiparty chat: Conversational agents in group settings with humans and models. *CoRR*, abs/2304.13835.

Thomas Wolf, Lysandre Debut, Victor Sanh, Julien Chaumond, Clement Delangue, Anthony Moi, Pierric Cistac, Tim Rault, Rémi Louf, Morgan Funtowicz, Joe Davison, Sam Shleifer, Patrick von Platen, Clara Ma, Yacine Jernite, Julien Plu, Canwen Xu, Teven Le Scao, Sylvain Gugger, Mariama Drame, Quentin Lhoest, and Alexander M. Rush. 2020. Transformers: State-of-the-art natural language processing. In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing: System Demonstrations, EMNLP 2020 - Demos, Online, November 16-20, 2020*, pages 38–45. Association for Computational Linguistics.

## A  Prompt

Here, we present the prompt that both the tested LLMs and the human participants received to complete the NSP task:

```
Your task is to predict the next
speaker given the full conversation
history. Do not provide any explanation.
Do not complete the conversation.

This is the conversation history:

<conversation history>

Predict the next speaker by outputting
the name and only the name of the
next speaker.  Carefully consider the
motives of the participating speakers
in the conversation.  Do not provide
any explanation.  Do not complete the
conversation.
```

## B  Inference and Training Details

- Hugging Face Transformers library (Wolf et al., 2020) for loading and running the models.[1]

- Inference
  - All models were loaded in 8-bit precision via bitsandbytes. [2]
  - Temperature: 0.0 (no sampling)

- Training
  - Supervised Fine-tuning Trainer script from Hugging Face Transformer Reinforcement Learning library. [3]
  - LoRA (Hu et al., 2022) with rank $r = 8$.

- Hardware equipment: Up to two NVIDIA RTX 6000 Ada Generation GPUs at the same time.

## C  Data Collection for Human Baseline

We describe our process of collecting data for the human baseline in detail.

---

[1] https://github.com/huggingface/transformers
[2] https://github.com/bitsandbytes-foundation/bitsandbytes
[3] https://github.com/huggingface/trl/

### C.1  Sample Selection

We targeted a sample of 1% of each test sets to keep the amount of work for the voluntary annotators small while still capturing the nature of the datasets. However, the different natures added additional constraints. For ML, we only selected three full dialogues leading to approximately 1% of the data. For AMI, a sample of 1% would have been outside of our annotator budget. Therefore, we selected a sample of 0.5%. For DiPCo, 1% was not enough to capture the dataset's nature, so we doubled the sample size here.

To decide which samples of the test sets to use during data collection, we performed several random samples of consecutive dialogue utterances and selected the one showing the most similar accuracy in a zero-shot setting to the full dataset.

### C.2  Annotation Acquisition

We asked colleagues working in the field of NLP and Computer Vision to fill out the forms for our user study to acquire a human baseline. The participation was not mandatory, and we offered no compensation. We informed the participants that the data created by them during this user study will be incorporated into a scientific publication.

We presented the participants of our data collection for the human baseline the following introduction texts:

- **Human Baseline for Next Speaker Prediction on the AMI Meeting Corpus Dataset**
  I want to acquire a human baseline for the task of next speaker prediction on (a sample of) the AMI Meeting Corpus Dataset (https://groups.inf.ed.ac.uk/ami/corpus/). You will be presented the same prompt as the models I am testing. You will then be asked to select the next speaker out of the given ones. Please read the instructions in the first prompt carefully. The following questions (63 in total) will have the same prompt and will only change the newest (and oldest) conversation step.

- **Human Baseline for Next Speaker Prediction on the Dinner Party Corpus (DiPCo) Dataset**
  I want to acquire a human baseline for the task of next speaker prediction on (a sample of) the Dinner Party Corpus (DiPCo) Dataset (https://arxiv.org/abs/1909.13447). You will

be presented the same prompt as the models I am testing. You will then be asked to select the next speaker out of the given ones. Please read the instructions in the first prompt carefully. The following questions (55 in total) will have the same prompt and will only change the newest (and oldest) conversation step.

- **Human Baseline for Next Speaker Prediction on the MultiLIGHT Dataset 1/3**
  I want to acquire a human baseline for the task of next speaker prediction on (a sample of) the MultiLIGHT dataset (https://arxiv.org/abs/2304.13835). This is the first of three full conversations I want your help for. You will be presented the same prompt as the models I am testing. You will then be asked to select the next speaker out of the given ones. Please only enter your choice after your prediction is final (as picking a choice will show the next conversation state). Please read the instructions in the first prompt carefully (whose conversation history will be empty). The following questions (26 in total) will only show the newest conversation step. You can use the conversation steps of the previous questions for your decision for the current question (as the models also get the full conversation history for the current conversation state).

- **Human Baseline for Next Speaker Prediction on the MultiLIGHT Dataset 2/3**
  I want to acquire a human baseline for the task of next speaker prediction on (a sample of) the MultiLIGHT dataset (https://arxiv.org/abs/2304.13835). This is the second of three full conversations I want your help for. You will be presented the same prompt as the models I am testing. You will then be asked to select the next speaker out of the given ones. Please only enter your choice after your prediction is final (as picking a choice will show the next conversation state). Please read the instructions in the first prompt carefully (whose conversation history will be empty). The following questions (31 in total) will only show the newest conversation step. You can use the conversation steps of the previous questions for your decision for the current question (as the models also get the full conversation history for the current conversation state).

- **Human Baseline for Next Speaker Prediction on the MultiLIGHT Dataset 3/3**
  I want to acquire a human baseline for the task of next speaker prediction on (a sample of) the MultiLIGHT dataset (https://arxiv.org/abs/2304.13835). This is the third of three full conversations I want your help for. You will be presented the same prompt as the models I am testing. You will then be asked to select the next speaker out of the given ones. Please only enter your choice after your prediction is final (as picking a choice will show the next conversation state). Please read the instructions in the first prompt carefully (whose conversation history will be empty). The following questions (34 in total) will only show the newest conversation step. You can use the conversation steps of the previous questions for your decision for the current question (as the models also get the full conversation history for the current conversation state).

The introduction text for ML differs from the other datasets as we used a different setup for the online form. This switch from the setup for ML to the one used for AMI and DiPCo was mostly done out of convenience during the creation of the online form and should not impact the results of the data collection.

After this introduction text, the participants were shown the exact same prompt template as they were presented to the LLMs (subsection 2.2). To select the next speaker, they could choose from all appearing speakers in that dialogue with a radio button control element.

# Adapting ASR Models for Speech-to-Punctuated-Text Recognition with Utterance Gluing

**Agata Jakubiak, Piotr Stachyra, Piotr Czubowski, Hubert Borkowski, Sebastian Łątka, Radosław Iżak**, **Kornel Jankowski**, **Sonia Janicka** and **Mateusz Zieliński**

Samsung R&D Institute Poland

{a.jakubiak2,p.stachyra,h.borkowski,s.latka}@partner.samsung.com,
{p.czubowski,r.izak,k.jankowski,s.janicka,m.zielinski3}@samsung.com

## Abstract

Punctuation prediction is a necessary part of ASR models, usually accomplished in a cascaded framework, where a secondary text-based model supplements an unpunctuated ASR output with punctuation marks. However, this approach results in ignoring acoustic context, which makes it poorly suited to certain languages. In this paper, we explore previously proposed ideas on an alternative approach, i.e. Speech-To-Punctuated-Text (STPT) models, and present a solution that allows adapting existing ASR models to output punctuated text. Additionally, we propose utterance gluing, a method of augmenting data to circumvent the lack of speech corpora with long utterances and punctuated references. Our STPT models trained on augmented data outperform STPT models trained on regular data, as well as traditional cascaded models, suggesting that acoustic-based punctuation prediction may be a good alternative to the more common text-based punctuation prediction.

## 1 Introduction

With the advances in Automatic Speech Recognition (ASR), speech recognition models have become useful in many contexts. Still, there are areas in ASR research which, despite their influence on practical usage, remain under-researched. One of these is punctuation prediction – the task of giving proper punctuation to the ASR output.

Appropriate punctuation in a text is important both for its readability to humans (Ákos Tündik et al., 2018), and for the success of downstream tasks which use it as input, such as machine translation (Vandeghinste et al., 2018) or named entity recognition (Nguyen et al., 2020). Long blocks of text, if not separated into sentences, can be difficult for humans and machines to parse through and understand; additionally, some sentences may be ambiguous without appropriate punctuation. For these reasons, no matter the use-case of an ASR

model, having a properly punctuated output is generally preferable.

Despite this, a still widely-used approach to ASR models is to make them output unpunctuated, lowercase text. Such text is often subject to a separate process called punctuation prediction (Gravano et al., 2009), which adds punctuation to it. Many punctuation prediction models do not use any acoustic features present in speech, relying only on the text output of ASR as their input; this is referred to as lexical punctuation prediction. However, this approach presents issues.

Firstly, if a text may be correctly punctuated in multiple ways, it is impossible for the model to distinguish between them without access to acoustic context. This is especially striking in languages that rely more heavily on the acoustic context rather than the grammatical structure of the sentence to disambiguate between different meanings, such as Spanish (Hualde, 2005) or French (Price, 2005), wherein questions are often distinguished from declarative statements exclusively through prosody.

Secondly, since the lexical punctuation prediction relies on the text output of the ASR, any ASR errors are likely to result in punctuation errors, as the punctuation prediction model tries to punctuate the incorrect sentence.

Thirdly, this approach adds the burden of maintaining an additional model alongside the ASR model itself. This is additionally problematic when working with limited memory and computational power, such as when running on mobile devices.

A practiced solution to the *first* and *second* issue is creating hybrid punctuation prediction models which use acoustic features as input alongside text (Klejch et al., 2017), and have access to additional acoustic context not present in the text itself. These models are usually bigger and more complex than purely lexical models, which makes the third issue even more prevalent. A less common solution, which addresses all three issues, is creating

ASR models that directly output punctuated text, and learn to place punctuation marks based solely on the speaker's prosody (Nozaki et al., 2022; Kim et al., 2023). This is referred to as acoustic punctuation prediction, and is the solution we are developing.

The biggest roadblock in developing robust Speech-To-Punctuated-Text (STPT) models is the lack of appropriate speech corpora with both punctuated references and long utterances. Discarding corpora without punctuation marks (e.g., LibriSpeech (Panayotov et al., 2015) and Multilingual LibriSpeech(Pratap et al., 2020)) means severely limiting training data, which unavoidably results in worse recognition metrics, especially in low-resource languages. Moreover, many widely-used speech corpora used for ASR training contain mostly one-sentence utterances (e.g., Common Voice (Ardila et al., 2019)). An STPT model trained on such a dataset is likely to learn to output periods and question marks at the ends of utterances only. This is usually not preferable, as most ASR models are unlikely to process only one sentence at a time.

In this paper, we propose a method of training an STPT model aimed at tackling both these issues without compromising on the Word Error Rate (WER) of the model.

## 2 Related Work

Creating an end-to-end ASR model that takes speech as input and outputs punctuated text has been previously undertaken for English and Japanese (Nozaki et al., 2022) and for English (Kim et al., 2023).

Mimura et al. (2021) tackled a close topic; however, their goals were much broader, including removal of filler words and changing the speech to be more formal, so their findings are largely inapplicable to our research.

Recently, STPT models have become much more popular, with models such as NVIDIA's Parakeet[1] and Canary[2] being published. These projects did not focus on punctuation; they used punctuated and capitalized transcripts as the training data, so the models learned to produce punctuation in the output, but the creators do not claim to have used any specific methods to improve punctuation results,

and they do not share any metrics showing their punctuation performance. We will be focusing on the punctuation-oriented research of Nozaki et al. (2022) and Kim et al. (2023) in our analysis.

### 2.1 Architecture changes

The main innovation suggested by Nozaki et al. (2022) on creating an STPT model is the addition of an auxiliary loss in an intermediate layer. In their experiments, this addition improved the performance of the model in multiple metrics; however, in the experiments conducted by Kim et al. (2023), the auxiliary loss did not seem to improve the performance of the model significantly.

Kim et al. (2023) focused on streaming, chunk-based ASR, in which their model was only provided with fragments of sentences at a time. This, as explored in more detail in Section 2.2, seems to make punctuation detection much more difficult.

### 2.2 Punctuation in long utterances

Nozaki et al. (2022) acknowledge that the English training corpus they use, MuST-C (Di Gangi et al., 2019), contains only single-sentence utterances, but they do not attempt to solve this issue. Their model achieves good results on single-sentence test cases, but they do not test it on longer utterances. Their Japanese test utterances are single-sentence only, while only one-sixth of the training ones contain more than one sentence.

Kim et al. (2023) also used MuST-C, but addressed the problem in two ways. Firstly, they concatenated random pairs of training utterances, so that every new utterance consisted of two sentences. Additionally, they also tested the model on long-form speech. The results on long-form test cases were worse than those achieved by Nozaki et al. (2022) on single-sentence test cases, particularly on periods and question marks. However, the model presented by Kim et al. (2023) achieved worse results on periods in single-sentence test cases than it did on periods in long-form test cases, which counter-intuitively suggests that it was actually better at predicting mid-utterance periods than it was at predicting utterance-ending ones. This is likely caused by the fact that its streaming ASR had access to less context, which made it difficult for the model to detect ends of utterances.

---

## 3 Proposed Method

Broadly speaking, we wanted our method to be as easy to adapt and use as possible. Because of that, the ideas we propose are focused on data processing, and could be implemented to add punctuation prediction to any ASR model; although, as mentioned before in relation to (Kim et al., 2023), some architectures seem better suited to the task of punctuation prediction than others.

### 3.1 Punctuation adaptation

In our research, we decided to adapt regular ASR models on punctuated data, rather than training STPT models from scratch. This has many advantages; namely, adapting a model for punctuation prediction is much faster and less resource-intensive than training an STPT model, which is practical for production contexts where time needed to deploy a new model is a factor. Additionally, with this method, training corpora without proper punctuation can still be used in the early phases of training to improve the final ASR model. Finally, with punctuation adaptation, anyone can add punctuation prediction to their existing ASR model, without restarting the training process, which makes the method easier to test and use.

### 3.2 Utterance gluing

As previously described, since many ASR corpora contain only one sentence in each utterance, STPT models trained on them struggle with placing periods and question marks in places other than the ends of utterances. Concatenating pairs of utterances has been proposed as a solution (Kim et al., 2023); however, an STPT model trained on concatenated utterances could learn to recognize artifacts generated by concatenation (e.g., changes of speakers, loudness, or in the background noise), and place punctuation there. We expanded on the idea of concatenation to make the final utterances resemble natural long-form speech in the following ways:

- Only utterances recorded by the same speaker are concatenated.

- Utterances shorter than $1\,\mathrm{second}$ and very quiet utterances (with RMS amplitude lower than 0.01) are discarded.

- Every speaker's utterances are sorted by RMS amplitude, and concatenated with the ones next to them on the sorted list, so that the concatenated utterances have similar volumes.

- Groups of variable numbers of utterances are concatenated, so that the model does not learn to rely on the number of sentences in an utterance.

- A short cross-fade (randomly chosen between 8, 10 and 12 ms) is added between the utterances.

- Long periods of silence from the resulting utterance are cut out, by randomly choosing $duration$ between 0.6, 0.7, 0.8 and 0.9 seconds, and cutting out all parts of the recording that are quieter than 0.2% of the maximum amplitude of a given recording and longer than $duration$. A fragment of silence $n$ seconds long (where $n$ is a random length shorter than $duration$) is left behind, so that some silence remains.

We call this method *utterance gluing*, as it is more complex than simple concatenation. The script used can be found online[3].

### 3.3 Data processing

We decided to support recognizing periods, commas, question marks, inverted question marks (¿), exclamation marks, and inverted exclamation marks (¡). Our data processing pipeline for punctuation data was as follows:

1. All punctuation marks other than those supported were removed from the reference text. Additionally, all periods used in abbreviations and initials were removed.

2. Every occurrence of a supported punctuation mark was replaced by a tag, written as a separate word; those tags were also placed in the token vocabulary of the model.

## 4 Models

### 4.1 ASR

The ASR model used in this work is a conformer-transducer, a sequence-to-sequence model, which is a variation of an architecture derived from the RNN-transducer (Graves, 2012). Specifically, we

---

[3] https://github.com/SamsungLabs/adapting-asr-models-for-stpt-with-utterance-gluing

employ the first-pass model architecture as described in section 2 in (Park et al., 2023) without a feedback path from the joiner to the predictor. We refrained from using the second-pass portion of the architecture, focusing on the applicability of the proposed method to a single-pass streaming model. We release the code used for training on GitHub[3].

The concept relies on employing transcriber and predictor networks: the former operating on the acoustic features $\mathcal{X} \in \mathbb{R}^d$ derived from the audio signal, the latter on the utterance transcription encoding $\mathcal{Y}$, representing wordpieces.

The transcriber takes an input sequence of acoustic features and outputs a transcription vector. In this work, the transcriber is a stack of 16 conformer layers (Gulati et al., 2020) capturing the global, as well as local patterns by utilizing attention and convolution layers. To ensure optimal resources utilization, we used striding as a reduction technique applied to the acoustic features, prior to processing by the transcriber.

The predictor consists of two layers of an LSTM network. Its purpose is to learn to model an output sequence $g = (g_0, g_1, ..., g_U)$, where $U$ corresponds to the tokens' sequence length.

It is worth noting that the input sequence is the original tokens' sequence $y = (y_1, ..., y_U)$ with an encoded null output $\emptyset$, prepended to it. Therefore, at the input, we process an extended input vector $\hat{y} = (\emptyset, y_1, ..., y_U)$, as proposed by previous work (Graves, 2012). Utilization of a blank token enables teaching the model how to align speech, i.e. account for silent parts in utterances without malforming the transcribed speech sequence in temporal context.

These networks are jointly trained using a Joiner, integrating the information from both networks, with an objective function (commonly known as RNN-T Loss) defined as log posterior probability: $\mathcal{L} = -ln(y|x)$. Joiner adds the outputs of transcriber and predictor, which are further passed through activation layer and linear layers.

The ASR we trained had 30 million parameters. An overview of the architecture used for the ASR model used in this work is shown in Figure 1.

### 4.2 Lexical restoration

To evaluate our approach against lexical methods, we also trained and tested transformer-based token classification models. This was done due to the lack of appropriate open-source models for this study;



Figure 1: Transducer architecture used in this work.

the most appropriate being KREDOR's punctuate-all model[4], based on (Guhr et al., 2021), which does not support exclamation marks and inverted punctuation marks. For each language, an instance of XLM-RoBERTa-large (Conneau et al., 2019) was first fine-tuned on a mix of long- and short-form utterances with a 1:4 ratio, and then further trained on the former only. The needed datasets were accessed through the OPUS (Tiedemann, 2012) website and included ParaCrawl (Bañón et al., 2020), OpenSubtitles (Lison and Tiedemann, 2016), and EuroParl (Koehn, 2005) to balance formal and informal writing styles. For each dataset, short-form sentences were retrieved and cleaned (e.g., abbreviations were removed). Then, a random subsample was concatenated to form utterances 2-6 sentences long. In total, each model was trained on more than 16 M utterances per epoch, with training ending after 15 epochs, or if the average of all punctuation mark metrics plateaued for more than two epochs.

## 5 Experiments

### 5.1 Datasets used

We decided to run our experiments on German, Polish, and Spanish, as those languages represent three different language subgroups (Eberhard et al., 2024), and we suspected that different approaches to punctuation prediction might work best for different kinds of languages. Unfortunately, we could not train an English model with MuST-C and compare it to previous works on this subject, (Nozaki

---

[4]https://huggingface.co/kredor/punctuate-all

et al., 2022) and (Kim et al., 2023), since the dataset is not currently available[5].

### 5.1.1 Training and validation datasets

For punctuation training purposes, we searched for open-source datasets with well-punctuated references. We decided to use Common Voice 16.1 (Ardila et al., 2019) for Spanish and German, and Common Voice 13.0 with ParlaSpeech (Ljubešić et al., 2025) for Polish. 1% of the data was selected for validation. The number of utterances and punctuation marks in each dataset can be seen in Table 1.

For the purposes of our experiments, we created four versions of each training and validation dataset:

1. A non-glued, non-punctuated version, used to train a regular ASR model.

2. A non-glued, punctuated version, with most of the utterances only containing one sentence, later referred to as "single-sentence punctuated data" (*single*).

3. A concatenated, punctuated version, where utterances were randomly concatenated into groups of 2-3, resulting in 361k utterances in German, 230k in Polish and 591k in Spanish, and their references concatenated accordingly (*concat*).

4. A glued, punctuated version, where utterances were glued together into groups of 2-3, using the methodology described in section 3.2, resulting in 339k utterances in German, 199k in Polish and 549k in Spanish, and their references concatenated accordingly (*glued*).

Table 1: Number of utterances and punctuation marks in original non-augmented datasets.

| Language | Utts | . | , | ¿ | ? | ¡ | ! |
|---|---|---|---|---|---|---|---|
| German | 867k | 801k | 218k | 0 | 47k | 0 | 22k |
| Polish | 556k | 446k | 578k | 0 | 51k | 0 | 69k |
| Spanish | 1418k | 1418k | 508k | 5.7k | 5.7k | 4.5k | 8.8k[6] |

---

### 5.1.2 Evaluation datasets

We needed to use real multi-sentence utterances to evaluate the models on actual mid-utterance periods, question marks and exclamation marks. We decided to use Multilingual LibriSpeech (MLS), which contains many long utterances from audiobooks (Pratap et al., 2020). The released version of this dataset does not contain punctuation in its references, but we restored the punctuation using the original books' text. Then, for each language, we selected 1024 utterances which contained at least one question mark from the training subset of the corpus, and we manually modified the references to only contain the punctuation marks we were using (e.g., replacing semicolons with periods). We did not simply remove the unsupported punctuation marks, as we did in training data, because MLS contained much more of them than our training datasets. However, we removed a few utterances which contained punctuation that could not be straightforwardly replaced. The dataset details can be seen in Table 2. The evaluation datasets were released on GitHub[3].

Table 2: Number of punctuation marks in evaluation datasets.

| Language | Utts | . | , | ¿ | ? | ¡ | ! |
|---|---|---|---|---|---|---|---|
| German | 1020 | 1825 | 3210 | 0 | 1421 | 0 | 429 |
| Polish | 1014 | 2958 | 4051 | 0 | 1364 | 0 | 351 |
| Spanish | 1022 | 2525 | 3134 | 1338 | 1338 | 323 | 323 |

## 5.2 Experiment methodology

In our experiment, we wanted to compare the effectiveness of the following approaches: *lexical* restoration and three variants of acoustic recognition: trained on *single*, *concat*, and *glued* punctuated data.

### 5.2.1 Acoustic model training

To that end, firstly, we trained a multilingual ASR model from scratch for 925k steps on the non-punctuated version of all three training datasets. Then, we adapted it on the non-punctuated training dataset for every language, resulting in three regular, non-punctuated ASR models. Then we adapted each of them on the *single*, *concat*, and *glued* punctuated data, resulting in three different STPT models for every language. Table 3 shows the numbers of training steps for each checkpoint chosen for evaluation.

### 5.2.2 Vocabulary

The token vocabulary of all of the models was the same. Tags used for punctuation prediction were present in the vocabulary from the start, and went unused by the earlier, non-punctuated models. Therefore, adaptations consisted simply of running training from a previously trained checkpoint, with entirely new training and validation data, and no other changes. When adapting a previously trained ASR model with no punctuation tags in the vocabulary, one could accomplish the same outcome by replacing the least used tokens in the vocabulary with punctuation tags. This would allow the model to adapt for punctuation prediction without the size of the vocabulary being changed, and without the need to retrain the model from scratch.

### 5.2.3 Lexical models

Additionally, for every language, we used our lexical punctuation prediction model (as described in Section 4.2) and KREDOR's punctuate-all model to create two cascaded, lexical STPT models out of the non-punctuated ASR models created in 5.2.1, in order to compare the acoustic models with state-of-the-art lexical punctuation prediction. It is worth mentioning that our lexical models are more than 18 times larger, and KREDOR is about 9 times larger, than our STPT models.

### 5.2.4 Performance metrics

To compare these approaches, we treated them as if the models were binary classifiers deciding whether or not the given punctuation mark should be placed at a given position in the recognized text and compared their precision, recall, and F1 scores. Additionally, we compared WERs of the models with punctuation marks excluded.

Table 3: Number of training steps for chosen checkpoints.

| Language | non-punct | single | concat | glued |
|---|---|---|---|---|
| German | 1891k | 2143k | 2000k | 1980k |
| Polish | 1569k | 1703k | 1600k | 1654k |
| Spanish | 1960k | 2420k | 2140k | 2155k |

### 5.3 Results and discussion

The evaluation results of the five previously described approaches for each language can be seen in Table 4. Since the lexical models used the outputs of non-punctuated ASR models, the WERs listed in the lexical models' rows are the WERs of acoustic models before punctuation adaptations.

They can be also used to see how punctuation adaptations affected WERs.

### 5.3.1 Exclamation marks

In our experiments, exclamation marks could not be reliably recognized by any model (best F1 score was 0.21, and most were far worse). In acoustic models, this does not seem to stem from them being underrepresented in training data (see Table 1). It is likely they are close enough to periods, both in their pronunciation and their usage, that neither lexical nor acoustic model can tell them apart. Since mistaking exclamation marks for periods does not usually impact the meaning of the text, we decided to treat exclamation marks as equivalent to periods in our results, and disregard inverted exclamation marks.

### 5.3.2 Lexical models

Our *lexical* models achieved similar results to KREDOR's state-of-the-art model, with the notable exception of question marks, where their results were better. For that reason, going forward, we will be using them as the lexical state-of-the-art benchmark. Although our models were trained on very similar data to each other, some metrics differ strongly between languages. This suggests that lexical punctuation prediction may be better suited for some languages than for others.

### 5.3.3 Acoustic models

In general, the *single* acoustic models performed very poorly, achieving the lowest F1 scores out of the acoustic models on all languages and punctuation marks, except for Spanish utterance-ending periods. As predicted, they were almost unable to produce mid-utterance periods and question marks, with the notable exception of Spanish mid-utterance periods.

In Polish and German, the *glued* models achieved the highest F1 scores on all punctuation marks, outperforming all other models, both acoustic and lexical. The most notable difference between *lexical* and *glued* models was in mid-utterance periods and mid-utterance question marks, though in Polish the difference on utterance-ending question marks was also large.

In Spanish, there is no clear best-performing model. Our Spanish acoustic models were by far the worst of the three languages at recognizing question marks, and they were outperformed by the *lexical* model. This is likely caused by question

Table 4: Comparison of recalls, precisions and F1 scores of punctuation marks' recognition between models. For sentence-ending punctuation marks, results are split into mid-utterance and utterance-ending marks. Exclamation marks have been treated as periods, and inverted exclamation marks have been deleted. WER values are calculated with punctuation marks excluded.

| Language | Model | WER | mid . | | | end . | | | , | | | mid ? | | | end ? | | | ¿ | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | | | Rec | Pre | F1 | Rec | Pre | F1 | Rec | Pre | F1 | Rec | Pre | F1 | Rec | Pre | F1 | Rec | Pre | F1 |
| German | KREDOR | 0.24 | 0.49 | 0.59 | 0.53 | 0.91 | 0.69 | 0.79 | **0.65** | 0.65 | **0.65** | 0.32 | 0.67 | 0.44 | 0.48 | 0.83 | 0.61 | - | - | - |
| | lexical | 0.24 | 0.50 | 0.53 | 0.52 | 0.88 | **0.72** | 0.79 | 0.64 | 0.67 | **0.65** | 0.41 | 0.62 | 0.49 | **0.53** | **0.88** | **0.66** | - | - | - |
| | single | 0.21 | 0.02 | **0.90** | 0.03 | 0.90 | 0.64 | 0.75 | 0.64 | 0.51 | 0.57 | 0.01 | 0.73 | 0.02 | 0.37 | 0.76 | 0.49 | - | - | - |
| | concat | **0.19** | 0.57 | 0.72 | 0.63 | **0.93** | 0.69 | 0.80 | 0.63 | 0.66 | 0.64 | 0.34 | **0.77** | 0.47 | 0.48 | 0.86 | 0.61 | - | - | - |
| | glued | **0.19** | **0.70** | 0.67 | **0.68** | **0.93** | 0.71 | **0.81** | 0.62 | **0.68** | 0.65 | **0.49** | 0.76 | **0.59** | 0.53 | 0.86 | **0.66** | - | - | - |
| Polish | KREDOR | 0.28 | 0.46 | 0.59 | 0.52 | 0.95 | 0.78 | 0.86 | 0.63 | 0.63 | 0.63 | 0.28 | 0.66 | 0.39 | 0.46 | 0.83 | 0.60 | - | - | - |
| | lexical | 0.28 | 0.46 | 0.57 | 0.51 | 0.94 | 0.78 | 0.85 | 0.61 | 0.62 | 0.62 | 0.39 | 0.61 | 0.47 | 0.48 | 0.85 | 0.61 | - | - | - |
| | single | 0.24 | 0.00 | 0.32 | 0.01 | 0.92 | 0.73 | 0.81 | 0.67 | 0.49 | 0.56 | 0.05 | **0.89** | 0.09 | 0.32 | 0.67 | 0.44 | - | - | - |
| | concat | 0.22 | 0.32 | **0.78** | 0.45 | 0.94 | 0.79 | 0.86 | **0.68** | 0.57 | 0.62 | 0.46 | 0.85 | 0.60 | 0.50 | 0.81 | 0.62 | - | - | - |
| | glued | **0.21** | 0.50 | **0.78** | **0.61** | **0.96** | **0.85** | **0.90** | 0.61 | **0.67** | **0.64** | **0.67** | 0.82 | **0.74** | **0.66** | **0.88** | **0.76** | - | - | - |
| Spanish | KREDOR | 0.24 | 0.39 | 0.55 | 0.45 | 0.99 | 0.54 | 0.70 | **0.52** | 0.47 | 0.50 | 0.06 | **0.59** | 0.11 | 0.07 | 0.88 | 0.14 | - | - | - |
| | lexical | 0.24 | 0.45 | 0.54 | 0.49 | 0.91 | 0.63 | 0.74 | 0.44 | 0.52 | 0.48 | **0.24** | 0.54 | **0.33** | **0.34** | 0.87 | **0.49** | 0.31 | 0.73 | **0.43** |
| | single | 0.33 | 0.27 | 0.68 | 0.39 | **0.99** | 0.64 | **0.78** | 0.36 | 0.50 | 0.42 | 0.02 | 0.41 | 0.03 | 0.01 | 0.75 | 0.02 | 0.03 | 0.62 | 0.05 |
| | concat | 0.17 | 0.52 | **0.76** | 0.62 | 0.97 | 0.58 | 0.72 | 0.51 | 0.54 | **0.53** | 0.20 | 0.44 | 0.28 | 0.21 | 0.87 | 0.34 | 0.25 | 0.62 | 0.36 |
| | glued | **0.16** | **0.74** | 0.63 | **0.68** | 0.98 | 0.56 | 0.71 | 0.40 | **0.60** | 0.48 | 0.22 | 0.55 | 0.32 | 0.14 | **0.88** | 0.24 | 0.24 | **0.75** | 0.36 |

marks being underrepresented in the Spanish training corpus. In internal experiments which utilized glued non-public data of better balance, higher results were achieved (0.39 recall and 0.88 precision for mid-utterance question marks, 0.38 recall and 0.94 precision for utterance-ending question marks, 0.35 recall and 0.88 precision for inverted question marks; for other punctuation marks, the results were comparable to the *glued* model).

### 5.3.4 Effects on WER

The WER seems positively affected by concatenation and gluing, although all acoustic models had access to the same training data, just processed differently. We think this is linked to the fact that the evaluation data consists of long utterances; it seems that training ASR models on long utterances improves their performance in recognizing long utterances.

### 5.3.5 Checkpoint instability

It is important to mention that during our training runs, the punctuation results between even close checkpoints varied strongly; it seemed difficult for an STPT model to find a local minimum for a punctuation task, as the model was trained for minimizing WER in general, without any special optimization for punctuation. It is likely that a training method with two loss functions, one aimed at minimizing WER and the other at optimizing the punctuation performance, could be used to improved the results further. That being said, we have trained our models for a significant time, and the

checkpoints we are presenting are the best of many, so we are reasonably sure that these are the best punctuation results possible with this method, despite the variability.

### 5.4 Possible new issues

We have found that acoustic punctuation prediction addresses issues inherent to lexical punctuation prediction, namely lexical ambiguity and dependence on good ASR output for good results. In our hands-on experiments, for example, a strong questioning tone of voice was enough to produce a question mark, regardless of whether the phrase spoken was grammatically a question, a statement, or even incoherent babble.

However, this approach creates new issues that need to be discussed. Some speakers may have a flat tone of voice that does not indicate a question when they are asking one. Some may pause while speaking, without intending for a comma or a period to be placed. In general, the performance of acoustic punctuation prediction is more dependent on the speaker, and how clearly they are speaking, and less dependent on whether the phrases they are using are grammatically correct, and have been recognized correctly.

Since we have proven that acoustic models can outperform lexical models, it seems that these issues are less prevalent than the ones present in lexical models, at least in our test cases.

## 6  Conclusions

In this paper, we postulate that acoustic punctuation prediction is a strong alternative to lexical punctuation prediction. We show that multi-sentence training utterances are necessary for training well-functioning STPT models, and that punctuated training corpora with single-sentence utterances can be augmented to be used for STPT model training. We theorize about the problems caused by concatenation, and we address them by developing our gluing technique. We show that gluing improves the results over concatenation (weighted avg F1 equal 0.5725 and 0.5371, respectively), and that both methods are superior to training acoustic models on single-sentence utterances. We also show that acoustic models can outperform lexical punctuation prediction models (with weighted avg F1 equal 0.4857), despite being much smaller.

## 7  Future work

The biggest challenge of end-to-end STPT models is the lack of well-punctuated corpora with multi-sentence utterances. This work was an attempt to circumvent that, and could be developed by improving the gluing methods further; however, if real long-utterance corpora were developed, the models trained on them would likely outperform the ones presented here, and possibly any model trained on glued data. Additionally, as we showed that languages can be better or worse suited for different approaches to punctuation prediction, we hope that more research on the topic will be conducted with non-English languages in mind.

Since the acoustic punctuation prediction is gaining popularity, as seen in models such as NVIDIA's Parakeet[1] and Canary[2], we believe it is important to measure and share the punctuation results of STPT models and work to improve these results, instead of treating punctuation as an afterthought. Judging by the high-quality outputs of these models, even though the authors did not share punctuation metrics, it seems that English STPT models can be trained on non-augmented punctuated data from scratch, since there is quite a large amount of such English data. For other languages, methods presented in this paper may be needed.

Lastly, we suggest that future efforts in developing speech corpora include punctuation in their references if possible, to enable further developments in this field.

## 8  Limitations

In our work, we have shown the advantage of acoustic models over lexical models when it comes to small ASR models trained on relatively small corpora, with relatively high WER. However, high WER negatively impacts the performance of lexical models, as the input they receive is unreliable. It would be useful to test these methods on larger, better-performing ASR models, and find if acoustic models continue to outperform lexical ones when the WER is lower.

Additionally, we have focused on one specific architecture – the sequence transducer – in our work. We hope the methods shown here are transferrable to different architectures, as none of our methods were reliant on the features of the sequence transducer. However, it is possible that different architectures differ in their suitability for use for STPT, and we do not know if the results shown here are representative of how every architecture would perform. This has to be investigated further to reach any definite conclusions.

## 9  Acknowledgements

## References

Rosana Ardila, Megan Branson, Kelly Davis, Michael Henretty, Michael Kohler, Josh Meyer, Reuben Morais, Lindsay Saunders, Francis M. Tyers, and Gregor Weber. 2019. Common voice: A massively-multilingual speech corpus. *ArXiv*, abs/1912.06670.

Marta Bañón, Pinzhen Chen, Barry Haddow, Kenneth Heafield, Hieu Hoang, Miquel Esplà-Gomis, Mikel L. Forcada, Amir Kamran, Faheem Kirefu, Philipp Koehn, Sergio Ortiz Rojas, Leopoldo Pla Sempere, Gema Ramírez-Sánchez, Elsa Sarrías, Marek Strelec, Brian Thompson, William Waites, Dion Wiggins, and Jaume Zaragoza. 2020. ParaCrawl: Web-scale acquisition of parallel corpora. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 4555–4567, Online. Association for Computational Linguistics.

---

[7]https://www.pytorchlightning.ai

Alexis Conneau, Kartikay Khandelwal, Naman Goyal, Vishrav Chaudhary, Guillaume Wenzek, Francisco Guzmán, Edouard Grave, Myle Ott, Luke Zettlemoyer, and Veselin Stoyanov. 2019. Unsupervised cross-lingual representation learning at scale. *arXiv*, abs/1911.02116.

Mattia A. Di Gangi, Roldano Cattoni, Luisa Bentivogli, Matteo Negri, and Marco Turchi. 2019. MuST-C: a Multilingual Speech Translation Corpus. In *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers)*, pages 2012–2017, Minneapolis, Minnesota. Association for Computational Linguistics.

David M Eberhard, Gary F Simons, and Charles D Fennig. 2024. Ethnologue: Languages of the World. Twenty-seventh edition. Dallas, Texas: SIL international.

Agustin Gravano, Martin Jansche, and Michiel Bacchiani. 2009. Restoring punctuation and capitalization in transcribed speech. *2009 IEEE International Conference on Acoustics, Speech and Signal Processing*, pages 4741–4744.

Alex Graves. 2012. Sequence transduction with recurrent neural networks. *ArXiv*, abs/1211.3711.

Oliver Guhr, Anne-Kathrin Schumann, Frank Bahrmann, and Hans Joachim Böhme. 2021. Fullstop: Multilingual deep models for punctuation prediction. In *Proceedings of the Swiss Text Analytics Conference 2021*, Winterthur, Switzerland. CEUR Workshop Proceedings.

Anmol Gulati, James Qin, Chung-Cheng Chiu, Niki Parmar, Yu Zhang, Jiahui Yu, Wei Han, Shibo Wang, Zhengdong Zhang, Yonghui Wu, and Ruoming Pang. 2020. Conformer: Convolution-augmented transformer for speech recognition.

José Ignacio Hualde. 2005. *The sounds of Spanish*. Cambridge University Press.

Jeff Hwang, Moto Hira, Caroline Chen, Xiaohui Zhang, Zhaoheng Ni, Guangzhi Sun, Pingchuan Ma, Ruizhe Huang, Vineel Pratap, Yuekai Zhang, Anurag Kumar, Chin-Yun Yu, Chuang Zhu, Chunxi Liu, Jacob Kahn, Mirco Ravanelli, Peng Sun, Shinji Watanabe, Yangyang Shi, and Yumeng Tao. 2023. TorchAudio 2.1: Advancing speech recognition, self-supervised learning, and audio processing components for Pytorch. In *2023 IEEE Automatic Speech Recognition and Understanding Workshop (ASRU)*, pages 1–9.

Hanbyul Kim, Seunghyun Seo, Lukas Lee, and Seolki Baek. 2023. Improved training for end-to-end streaming automatic speech recognition model with punctuation. In *Interspeech 2023*, pages 1653–1657.

Ondřej Klejch, Peter Bell, and Steve Renals. 2017. Sequence-to-sequence models for punctuated transcription combining lexical and acoustic features. In

*2017 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, pages 5700–5704.

Philipp Koehn. 2005. Europarl: A parallel corpus for statistical machine translation. In *Proceedings of Machine Translation Summit X: Papers*, pages 79–86, Phuket, Thailand.

Pierre Lison and Jörg Tiedemann. 2016. OpenSubtitles2016: Extracting large parallel corpora from movie and TV subtitles. In *Proceedings of the Tenth International Conference on Language Resources and Evaluation (LREC'16)*, pages 923–929, Portorož, Slovenia. European Language Resources Association (ELRA).

Nikola Ljubešić, Peter Rupnik, and Danijel Koržinek. 2025. The ParlaSpeech collection of automatically generated speech and text datasets from parliamentary proceedings. In *Speech and Computer*, pages 137–150, Cham. Springer Nature Switzerland.

Masato Mimura, Shinsuke Sakai, and Tatsuya Kawahara. 2021. An end-to-end model from speech to clean transcript for parliamentary meetings. In *2021 Asia-Pacific Signal and Information Processing Association Annual Summit and Conference (APSIPA ASC)*, pages 465–470.

Thai Binh Nguyen, Quang Minh Nguyen, Thi Thu Hien Nguyen, Quoc Truong Do, and Chi Mai Luong. 2020. Improving vietnamese named entity recognition from speech using word capitalization and punctuation recovery models. In *Interspeech 2020*, pages 4263–4267.

Jumon Nozaki, Tatsuya Kawahara, Kenkichi Ishizuka, and Taiichi Hashimoto. 2022. End-to-end speech-to-punctuated-text recognition. In *Interspeech 2022*, pages 1811–1815.

Vassil Panayotov, Guoguo Chen, Daniel Povey, and Sanjeev Khudanpur. 2015. Librispeech: an ASR corpus based on public domain audio books. In *2015 IEEE international conference on acoustics, speech and signal processing (ICASSP)*, pages 5206–5210. IEEE.

Jinhwan Park, Sichen Jin, Junmo Park, Sungsoo Kim, Dhairya Sandhyana, Changheon Lee, Myoungji Han, Jungin Lee, Seokyeong Jung, Changwoo Han, and Chanwoo Kim. 2023. Conformer-based on-device streaming speech recognition with kd compression and two-pass architecture. In *2022 IEEE Spoken Language Technology Workshop (SLT)*, pages 92–99.

Adam Paszke, Sam Gross, Soumith Chintala, Gregory Chanan, Edward Yang, Zachary DeVito, Zeming Lin, Alban Desmaison, Luca Antiga, and Adam Lerer. 2017. Automatic differentiation in PyTorch. In *NIPS-W*.

Vineel Pratap, Qiantong Xu, Anuroop Sriram, Gabriel Synnaeve, and Ronan Collobert. 2020. MLS: A large-scale multilingual dataset for speech research. In *Interspeech 2020*, pages 2757–2761.

Glanville Price. 2005. *Intonation*, chapter 20. John Wiley & Sons, Ltd.

Ole Tange. 2011. Gnu parallel - the command-line power tool. *login: The USENIX Magazine*, February 2011:42–47.

Jörg Tiedemann. 2012. Parallel data, tools and interfaces in OPUS. In *Proceedings of the Eighth International Conference on Language Resources and Evaluation (LREC'12)*, pages 2214–2218, Istanbul, Turkey. European Language Resources Association (ELRA).

Vincent Vandeghinste, Lyan Verwimp, Joris Pelemans, and Patrick Wambacq. 2018. A comparison of different punctuation prediction approaches in a translation context. In *Proceedings of the 21st Annual Conference of the European Association for Machine Translation*, pages 289–298, Alicante, Spain.

Máté Ákos Tündik, György Szaszák, Gábor Gosztolya, and András Beke. 2018. User-centric evaluation of automatic punctuation in asr closed captioning. In *Interspeech 2018*, pages 2628–2632.

# Word-level Language Identification using Encoder-Only and Decoder-Only Models

**Javier Iranzo-Sánchez, Parnia Bahar,**
**Alejandro Pérez-González-de-Martos**, **Mattia Antonino Di Gangi**
AppTek
{jiranzo,pbahar,aperez,mdigangi}@apptek.com

## Abstract

Language identification is a task that often finds applications in NLP pipelines that serve multiple languages. The task is classically presented as a sentence classification problem and models' performance degrades quickly when applying them to short phrases or individual words. Although challenging, fine-grained language identification is key to improve the performance of downstream tasks. This work explores the performance of both Encoder-Only and Decoder-Only Transformer Language models for the task of automatic word-level language identification. The results show that for this particular task, small Encoder-Only models outperform larger Decoder-Only models.

## 1   Introduction

This paper explores Word-Level Language Identification (WLID) within the context of a cascaded Speech-to-Speech (S2S) translation system with human supervision as an example application. Although there are several promising end-to-end approaches, the cascaded approach remains the preferred choice when human intervention is desired at multiple steps of the process. For the Speech-to-Speech or dubbing task, an additional problem occurs when the text to be uttered automatically contains words belonging to a language other than the target language. These words are a source of errors because the normal rules for pronunciation of the target language cannot be applied. There are many possible sources for these words, such as named entities, slang and loanwords. Fine-grained language labels can enhance various applications, including Text-to-Speech (TTS) models, by generating more accurate phoneme sequences (Vesik et al., 2020; Zhu et al., 2022) or using language-specific embeddings (Yang et al., 2024).

The contributions of this paper are three-fold: 1) We annotate a novel dataset for the word-level language identification task under the translation setting, 2) we benchmark multiple automatic approaches to this problem, including both Encoder-Only and Decoder-Only Large Language Models (LLMs) and 3) we propose new techniques to alleviate LLMs hallucinations in the context of the WLID task.

### 1.1   Related work

To the best of our knowledge, there are no works that address the WLID task in the context of dubbing. The closest related task is code-switching identification, which we take as a starting point since it is the most similar. There are however significant differences between the two. Code-switching is a stylistic choice of the speaker, typically used in informal contexts, whereas this work deals with the presence of foreign words within text in the target language, which mainly occurs as a result of the translation of foreign media. Code-switching techniques and models can thus be used for this task, but the difference in domains and formality levels means that the techniques and findings of the standard code-switching approaches might not translate to this specific task. This motivates the need for specific training and evaluation data to assess and improve the performance of automatic systems.

Automatic approaches to code-switching can include both hand-crafted rules and statistical models, as well as hybrid systems that combine the two. Iliescu et al. (2021) compare multiple approaches using semi-supervised data, whereas Osmelak and Wintner (2023) train a Conditional Random Field system whose input is a sequence of word-level features. Sterner and Teufel (2023) proposed a rule-based system (TongueSwitcher) and compared it with a BERT-like model trained on the data labeled with TongueSwitcher and human labels, and observed similar performance for German-English. Additionally, much work has been done to study the effects of code-switched text on the performance

Table 1: Dataset statistics, including number of sentences (*#sent*), number of words (*#word*s) and the number of those words that have been tagged as English (*#En words*).

| | Spanish - TED | | | Spanish - Media | | | German - TED | | |
|---|---|---|---|---|---|---|---|---|---|
| | #sent | #words | #En words | #sent | #words | #En words | #sent | #words | #En words |
| train | 2048 | 39182 | 1849 | - | - | - | 1024 | 17719 | 1014 |
| dev | 1316 | 26076 | 310 | - | - | - | 1574 | 25269 | 300 |
| test | 2502 | 42294 | 454 | 1854 | 9959 | 139 | 2823 | 43197 | 575 |

of automatic models. Winata et al. (2021) compare multiple techniques and finds that good results are obtained with the XLM-RoBERTa family of models. Zhang et al. (2023) find that LLM's performance significantly decreases for code-switched data across a variety of tasks (Sentiment Analysis, Machine Translation, Summarization and Code-Switching Language Identification). Their results show that it is competitive to finetune a smaller model rather than using an LLM. In the present work, we explore further the relative performance of Encoder-Only models and larger LLMs (Decoder-Only) using different approaches.

## 2 Methodology

### 2.1 Datasets

The main dataset used for the experiments reported on this paper is the MuST-C dataset (Di Gangi et al., 2019), a Speech Translation dataset that contains the recordings of multiple English TED talks as well as their translations into multiple languages. Specifically, we used the English-Spanish and English-German translation sets. We also experimented with an in-house dataset of media content. This dataset consists of English media with translations into Spanish.

The original MuST-C dataset does not include WLID labels, so we asked 2 native speakers of the target language to annotate each set. Table 1 reports a summary of the dataset statistics. The majority of the words are in Spanish, with around 1% of the words being in English. However, 10% of the sentences contain at least 1 English word, so even if the amount of words is low, it is common enough that the user-perceived quality is affected if this issue is neglected. The manually annotated training set was constructed so that there is a 1:1 proportion between sentences with and without English words. The remaining MuST-C train sentences were automatically annotated with Llama 3.1 70B, to be

used for semi-supervised experiments. [1]

### 2.2 Models

Both Encoder-Only and Decoder-Only models are tested based on previous results from the literature. For the first case, XLM-RoBERTa (Conneau et al., 2020) was used, in both *base* (270M) and *large* (550M) configurations. We take the pre-trained model and fine-tune it for the WLID task following a token classification approach, similarly to what is done for Named Entity Recognition (NER). Additionally, the existing Encoder-Only TongueSwitcher (Sterner and Teufel, 2023) model is also tested, which is a multilingual BERT model (Devlin et al., 2019) (172M) German-English code-switching model. The TongueSwitcher model has two versions: a pre-trained version that has been trained for the language modeling task with 24.6M Tweets that contain mixed German and English, and a code-switch detection model that has been further fine-tuned with supervised code-switching annotations. For the second case, we used Decoder-Only LLM from the Llama family. The recently released Llama3.1 (Dubey et al., 2024) 8B and 70B models were selected. After iterating through multiple prompts, we ended up with the prompt format shown in Table 2. Making the model output a label for every word in the sentence rather than only those on a different language, as well as forcing the output to be generated in a CSV-like format were significantly helpful to improve the accuracy of the model and to ensure that the model copies the input sentence.

Even after iterating multiple times to find the optimal prompt, we still observe many occurrences of hallucinations, that is, the generation of a sequence of words that differs from the original sentence to be annotated. This is not acceptable because the WLID system should add language annotations if

---

[1]The labels to reproduce the dataset are made available at https://github.com/mattiadg/wlid-annotations.

Table 2: Prompt format used for LLM inference.

| | |
|---|---|
| Instruction | The input is a {default_language} sentence. Your task is to output the language for each word in the sentence. Write one line for each word in the original sentence. Each output line will contain the word and the language, separated by a comma and a space. If a word exists in {default_language} and other languages, write {default_language}. Only answer to the last question and do not write additional questions. |
| Input | He comprado un ordenador ThinkPad. |
| Response | He, Spanish<br>comprado, Spanish<br>un, Spanish<br>ordenador, Spanish<br>ThinkPad., English. |

necessary, but leave the input text unchanged otherwise. We propose two techniques to post-process LLM hypothesis for which hallucinations are detected. The first is to replace the LLM hypothesis by the default hypothesis, which is the one where no words are labelled as a foreign language. As a second technique, we propose a post-processing algorithm called AutoMap to match the generated text against the original sentence. Specifically, we initially assign the default target language label to every word on the original sentence. Then, we take each generated word and compare it with the words in the original sentence. If there is a match, we assign the label of the generated word. Figure 1 provides an example of *AutoMap* in action.

## 3 Experiments

All development decisions are made based on the results on the MuST-C dev set. XLM-RoBERTa models are trained with Adam (Kingma and Ba, 2015) using 1e-5 learning rate and batch size 16, for a total of 8k steps with early-stopping every 500 steps. The learning rate is linearly scaled during the first 10% steps. Table 3 reports the results for the XLM-RoBERTa model based on the number of available training samples. Additionally, we also test wheter using the semi-supervised data annotated with Llama 3.1 is helpful, by adding 2048 sentences to the largest configuration, for a total of 4096 sentences (+*SSup*). Results are reported using the F1 score of the English class, as all of the tested configurations achieve 1.00 F1 score for the non-English class after rounding-up. The model is able to obtain acceptable results starting from 128 training samples, with increases in quality each time the available data doubles in size, starting to plateau when reaching 2048. Adding additional semi-supervised data degrades the performance rather than helping.

Table 3: XLM-RoBERTa results on the MuST-C Spanish dev set, using either the **B**ase or the **L**arge configuration. +*SSup* includes an additional 2048 examples automatically annotated with Llama. F1 scores for the English class.

| | Number of training samples | | | | | |
|---|---|---|---|---|---|---|
| | 128 | 256 | 512 | 1024 | 2048 | +SSup |
| **B** | 0.73 | 0.75 | 0.78 | 0.81 | 0.82 | 0.62 |
| **L** | 0.77 | 0.80 | 0.80 | 0.82 | 0.83 | 0.67 |

LLM models were tested both using the in-context learning (ICL) approach as well as fine-tuning (FT) with LoRA (Hu et al., 2022). Sampling is disabled when generating the LLM hypothesis, as we found that this helped to slightly increase quality and reduce hallucinations. Table 4 shows the performance of the LLM ICL approach on the MuST-C dev set. The train subset was shuffled once and then the first $n$ samples were selected to be used in the prompt. That is, the example selected for $n = 1$ is also used for $n = 2$ and so on. We observe no performance improvements for increasing the number of examples beyond 1.

Table 4: LLM evaluation results for MuST-C Spanish dev set, using $n$ in-context samples. Results show English-class F1 score.

| | $n$ | | | | | |
|---|---|---|---|---|---|---|
| Model | 1 | 2 | 4 | 8 | 16 | 32 |
| L-8B | 0.54 | 0.52 | 0.49 | 0.47 | 0.50 | 0.50 |
| L-70B | 0.71 | 0.71 | 0.70 | 0.70 | 0.70 | 0.69 |

For fine-tuning with LoRA, the best results were obtained with learning rate 1e-4, rank 16, $\alpha = 32$, dropout 0.05 and 8 epochs of fine-tuning. Table 5 compares the results of both ICL and FT depending on the post-processing technique. The results high-

Figure 1: Example of AutoMap post-processing for LLM hallucination. The labels of the LLM hypothesis (bottom) are mapped to the original text (top) by looking for exact matches (ignoring casing and punctuation) between the hypothesis and the original text. The text is in Spanish and the shaded box represents a word detected as English. The LLM hallucinated and failed to generate a label for *"un browser, prueba con"*, which also includes the English word *browser*, so it retains the default labels for those words.

Table 5: LLM performance on the MuST-C Spanish dev set. We compare scoring the raw output (∅), using AutoMap with exact matches (Ams) and using AutoMap but ignoring casing and punctuation (Am). F1 scores for the English class.

|  | | 8B | | | 70B | |
|---|---|---|---|---|---|---|
|  | ∅ | Ams | Am | ∅ | Ams | Am |
| ICL | 0.01 | 0.45 | 0.54 | 0.23 | 0.60 | 0.71 |
| FT | 0.01 | 0.45 | 0.59 | 0.01 | 0.56 | 0.72 |

Table 6: Final evaluation results on the test sets, for XLM-RoBERTa (R-Base, R-Large) and Llama3.1 (L-8B, L-70B) models. Precision/Recall for the English class.

|  | Spanish | | | | German | |
|  | Ted | | Media | | Ted | |
| Model | P | R | P | R | P | R |
|---|---|---|---|---|---|---|
| R-Base | 0.68 | 0.94 | 0.69 | 0.91 | 0.62 | 0.92 |
| R-Large | 0.69 | 0.98 | 0.73 | 0.94 | 0.68 | 0.92 |
| L-8B | 0.40 | 0.93 | 0.60 | 0.86 | 0.42 | 0.95 |
| L-70B | 0.48 | 0.97 | 0.68 | 0.86 | 0.45 | 0.96 |
| TS | - | - | - | - | 0.64 | 0.49 |
| FT-TS | - | - | - | - | 0.73 | 0.86 |

light the importance of the AutoMap technique in mitigating hallucinations. It can be observed how results are very poor without AutoMap, as the model struggles to reproduce the input sentence. However, the introduction of AutoMap (Ams) significantly boosts the performance of the system. Results are improved further if punctuation and casing are not taken into account when looking for word matches (Am), which indicates that casing and punctuation account for a significant portion of the mistakes. When using AutoMap, the finetuned models improve the ICL results by 0.05 F1 for the 8B model, and 0.01 F1 for the 70B model. Once again, this highlights the importance of AutoMap, as it allows to extract better performance from the fine-tuned models. The results also suggest that fine-tuning is able to increase the linguistic knowledge of the model, which helps to better detect foreign words, but it is not helpful for the model to learn to copy the input.

Table 6 shows the evaluation of the final models on the selected test sets. The English-German models are also compared with two versions of TongueSwitcher: the code-switch detection BERT-based model (TS) pre-trained on ample English-German code-switching data, as well as the baseline TS model fine-tuned with our WLID data (FT-

TS). Similarly to what was observed on the dev set, RoBERTa-based models outperform the Llama 3 models on the TED talks evaluation set, both for the Spanish and the German case. The TS code-switching system underperforms the other systems, and its performance only recovers when it has been trained with our WLID data (FT-TS). This highlights the need for specific data for WLID, as the existing code-switching systems cannot be directly applied to this task.

## 4  Conclusions

This work has introduced a new setting for word-level language identification, and provided a set of in-depth experiments to assess the performance of automatic models. Two interesting findings arise out of this research. First, there is still room for improvement on this task, on both the in-domain talks and out-of-domain media settings. Secondly, unlike current trends that tend to favor Decoder-Only LLMs, Encoder-Only models are a competitive, cost-efficient alternative for this task.

In terms of future work, Encoder-Only models

can be extended to the multilingual setting in order to simplify deployment, reduce costs and to improve quality and robustness. Additionally, the performance of both Encoder-Only and Decoder-Only models should be tested on a zero-shot setting, to assess their capabilities on language pairs for which little or no training data exists.

# References

Alexis Conneau, Kartikay Khandelwal, Naman Goyal, Vishrav Chaudhary, Guillaume Wenzek, Francisco Guzmán, Edouard Grave, Myle Ott, Luke Zettlemoyer, and Veselin Stoyanov. 2020. Unsupervised cross-lingual representation learning at scale. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 8440–8451, Online. Association for Computational Linguistics.

Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. 2019. BERT: Pre-training of deep bidirectional transformers for language understanding. In *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers)*, pages 4171–4186, Minneapolis, Minnesota. Association for Computational Linguistics.

Mattia A. Di Gangi, Roldano Cattoni, Luisa Bentivogli, Matteo Negri, and Marco Turchi. 2019. MuST-C: a Multilingual Speech Translation Corpus. In *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers)*, pages 2012–2017, Minneapolis, Minnesota. Association for Computational Linguistics.

Abhimanyu Dubey et al. 2024. The llama 3 herd of models. *arXiv preprint arXiv:2407.21783v1*.

Edward J. Hu, Yelong Shen, Phillip Wallis, Zeyuan Allen-Zhu, Yuanzhi Li, Shean Wang, Lu Wang, and Weizhu Chen. 2022. Lora: Low-rank adaptation of large language models. In *The Tenth International Conference on Learning Representations, ICLR 2022, Virtual Event, April 25-29, 2022*. OpenReview.net.

Dana-Maria Iliescu, Rasmus Grand, Sara Qirko, and Rob van der Goot. 2021. Much gracias: Semi-supervised code-switch detection for Spanish-English: How far can we get? In *Proceedings of the Fifth Workshop on Computational Approaches to Linguistic Code-Switching*, pages 65–71, Online. Association for Computational Linguistics.

Diederik P. Kingma and Jimmy Ba. 2015. Adam: A method for stochastic optimization. In *3rd International Conference on Learning Representations, ICLR 2015, San Diego, CA, USA, May 7-9, 2015, Conference Track Proceedings*.

Doreen Osmelak and Shuly Wintner. 2023. The denglisch corpus of German-English code-switching. In *Proceedings of the 5th Workshop on Research in Computational Linguistic Typology and Multilingual NLP*, pages 42–51, Dubrovnik, Croatia. Association for Computational Linguistics.

Igor Sterner and Simone Teufel. 2023. TongueSwitcher: Fine-grained identification of German-English code-switching. In *Proceedings of the 6th Workshop on Computational Approaches to Linguistic Code-Switching*, pages 1–13, Singapore. Association for Computational Linguistics.

Kaili Vesik, Muhammad Abdul-Mageed, and Miikka Silfverberg. 2020. One model to pronounce them all: Multilingual grapheme-to-phoneme conversion with a transformer ensemble. In *Proceedings of the 17th SIGMORPHON Workshop on Computational Research in Phonetics, Phonology, and Morphology*, pages 146–152, Online. Association for Computational Linguistics.

Genta Indra Winata, Samuel Cahyawijaya, Zihan Liu, Zhaojiang Lin, Andrea Madotto, and Pascale Fung. 2021. Are multilingual models effective in code-switching? In *Proceedings of the Fifth Workshop on Computational Approaches to Linguistic Code-Switching*, pages 142–153, Online. Association for Computational Linguistics.

Huai-Zhe Yang, Chia-Ping Chen, Shan-Yun He, and Cheng-Ruei Li. 2024. Bilingual and code-switching tts enhanced with denoising diffusion model and gan. In *Interspeech 2024*, pages 4938–4942.

Ruochen Zhang, Samuel Cahyawijaya, Jan Christian Blaise Cruz, Genta Winata, and Alham Aji. 2023. Multilingual large language models are not (yet) code-switchers. In *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing*, pages 12567–12582, Singapore. Association for Computational Linguistics.

Jian Zhu, Cong Zhang, and David Jurgens. 2022. ByT5 model for massively multilingual grapheme-to-phoneme conversion. In *Proc. Interspeech 2022*, pages 446–450.

# Beyond Shallow Heuristics: Leveraging Human Intuition for Curriculum Learning

**Vanessa Toborek[1,2], Sebastian Müller[1,2], Tim Selbach[1],**
**Tamás Horváth[1,2,3], Christian Bauckhage[1,2,3]**
[1]University of Bonn, [2]Lamarr Institute, [3]Fraunhofer IAIS
toborek@cs.uni-bonn.de

## Abstract

Curriculum learning (CL) aims to improve training by presenting data from "easy" to "hard", yet defining and measuring linguistic difficulty remains an open challenge. We investigate whether human-curated simple language can serve as an effective signal for CL. Using the article-level labels from the Simple Wikipedia corpus, we compare label-based curricula to competence-based strategies relying on shallow heuristics. Our experiments with a BERT-tiny model show that adding simple data alone yields no clear benefit. However, structuring it via a curriculum – especially when introduced first – consistently improves perplexity, particularly on simple language. In contrast, competence-based curricula lead to no consistent gains over random ordering, probably because they fail to effectively separate the two classes. Our results suggest that human intuition about linguistic difficulty can guide CL for language model pre-training.

## 1 Introduction

The growing scale of language models (LMs) has increased interest in training strategies that improve efficiency and convergence. Curriculum learning (CL), inspired by developmental psychology, is one such approach. CL structures training by presenting examples in a sensible order – typically from "easy" to "hard" (Elman, 1993; Bengio et al., 2009; Wang et al., 2021). While intuitively compelling and empirically useful in certain NLP tasks (Platanios et al., 2019; Nagatsuka et al., 2021), its overall impact on masked language model (MLM) pre-training remains debated (Surkov et al., 2022).

A key challenge in CL is the definition of linguistic difficulty. Unlike other domains, language difficulty may arise from multiple dimensions – such as syntax, semantics or context. In the absence of gold standards, prior work often relies on shallow heuristics (Platanios et al., 2019; Ranaldi

| Rarity | Class | Example |
|--------|-------|---------|
| low | SL | She is the author of the Twilight series. |
| low | EL | The history of poker is the subject of some debate. |
| high | SL | Today, most automotive diesels are turbocharged. |
| high | EL | Pink Floyd watched The Beatles recording Lovely Rita. |

Table 1: Sentences showing examples of high and low average word rarity for each class in the Simple Wikipedia dataset (Kauchak, 2013).

et al., 2023). Yet, readability research suggests that no single heuristic reliably captures linguistic complexity (Battisti et al., 2020). In contrast, humans intuitively consider multiple dimensions when simplifying text. This motivates the central question for this work: *Can human-curated simple language effectively guide CL for MLM pre-training?*

To answer this question, we study CL strategies based on article-level labels from the Simple Wikipedia corpus (Coster and Kauchak, 2011) and compare them to *competence-based CL* with shallow difficulty heuristics (Platanios et al., 2019), using BERT-tiny for MLM pre-training. Our experiments show that merely adding simple language data to training yields no overall improvement. Still, incorporating it through a label-based curriculum consistently improves not only overall perplexity but particularly the simple language perplexity. This effect vanishes when reversed: training on everyday language first is detrimental to learning, underscoring the importance of example ordering. Surprisingly, competence-based curricula show no benefit over random ordering.

Further, we find that simple and everyday language articles have similar vocabulary sizes and

high lexical and distributional overlap on the chosen difficulty heuristics. This suggests that competence-based CL fails here, because the heuristics do not effectively separate the classes. In contrast, the consistent gains from label-based curricula imply that simple language encodes other useful information, providing structure that benefits pre-training when leveraged correctly. These results suggest that simple language does indeed help, when applied in a curriculum that makes use of human intuition on linguistic difficulty.

## 2 Related Work

A common form of data-level CL orders the data points according to a global difficulty measure. This approach has been applied to various NLP tasks such as language modelling (Nagatsuka et al., 2021; Ranaldi et al., 2023), machine translation (Platanios et al., 2019; Mohiuddin et al., 2022), and questions answering (Liu et al., 2018) using difficulty measures like input length (Nagatsuka et al., 2021; Zaremba and Sutskever, 2015), word rarity (Platanios et al., 2019), or domain similarity (Mohiuddin et al., 2022). However, the choice of metric is often intuitive and its overall effectiveness remains debated, as the work by Surkov et al. (2022) found that competence-based CL for MLM offers little to no benefit.

A parallel line of work explores the benefits of simplified language in neural network training. Mueller and Linzen (2023) show that pre-training on simple language corpora strengthens the syntactic inductive bias in encoder-decoder models. Huebner et al. (2021) demonstrate that child-directed data facilitates grammar learning for down-sized encoder-only models. Lucas et al. (2024) explore CL through a masking-based strategy, also leveraging simplified language. While these studies focus on specific linguistic gains or efficiency improvements, the role of simplified language in global, data-level curriculum design remains unexplored. We address this gap by investigating whether editorially curated simple language – such as that in Simple Wikipedia – can serve as an effective learning signal for CL, and how it compares to commonly used difficulty heuristics.

## 3 Methodology

We use the following experimental setup to study the effect of simple language in MLM pre-training.

| Label | # tokens | # sentences |
|---|---|---|
| Simple (SL) | $3,395,297$ | $191,318$ |
| Everyday (EL) | $3,796,654$ | $176,019$ |

Table 2: Dataset statistics for simple (SL) and everyday (EL) language in the Simple Wikipedia corpus.

**Dataset** We employ the Simple Wikipedia dataset (Coster and Kauchak, 2011), the most popular, freely available simple language corpus in English. It consists of articles from the Simple English Wikipedia in simple language (SL) and their counterparts from the English Wikipedia in everyday language (EL). Each sentence inherits the article-level label (SL or EL), which may introduce some label noise due to within-article variation in sentence complexity. Table 2 compares both classes regarding their respective number of tokens and sentences.

**Difficulty Heuristics** For the competence-based CL, we consider three shallow heuristics for text difficulty: sentence length, word rarity, and the Flesch Reading Ease (FRE) score (cf. Platanios et al. (2019), Ranaldi et al. (2023)). Refer to Appendix B for the details. In addition to these, we include a random baseline, where difficulty scores are sampled uniformly to isolate the effect of data ordering from the progressive data exposure.

**Curriculum Strategies** We compare two CL paradigms. First, following Platanios et al. (2019), we implement the *competence-based* curriculum approach. We sort the training examples according to the aforementioned difficulty measures and gradually expand the training set as model competence increases. The curriculum proceeds until the entire dataset is included. We provide the full implementation details in Appendix A.

Second, we implement two *label-based* curricula using the SL/EL distinction. The sequential strategy first trains on SL until convergence, then continues training on EL. To mitigate potential forgetting from fully replacing the training data, we propose an incremental strategy: the model is first trained on SL alone, then continues on the combined SL+EL set, each phase until convergence. We also include a reverse sequential strategy (first on EL, then SL) as a control strategy.

**Training Setup** We train a BERT-tiny model with two transformer layers of hidden size 128,

| Strategy | Perplexity | SL Perplexity | EL Perplexity | # Updates |
|---|---|---|---|---|
| Baseline EL | 69.25 $_{\pm 4.04}$ | 59.50 $_{\pm 4.38}$ | 81.78 $_{\pm 4.85}$ | 658 667 $_{\pm 113\,192}$ |
| Baseline SL+EL | 69.61 $_{\pm 4.87}$ | 64.15 $_{\pm 5.05}$ | 76.46♪ $_{\pm 5.28}$ | 665 333 $_{\pm 102\,111}$ |
| Incremental | 66.36 $_{\pm 2.53}$ | 63.29 $_{\pm 3.39}$ | **71.51**♪ $_{\pm 2.55}$ | 781 333 $_{\pm 83\,312}$ |
| Sequential | **65.31**♪ $_{\pm 4.19}$ | **57.83**♪ $_{\pm 4.52}$ | 74.39 $_{\pm 4.91}$ | 781 333 $_{\pm 122\,292}$ |
| Anti-Sequential | 70.32 $_{\pm 3.97}$ | 59.24 $_{\pm 4.01}$ | 81.70♪ $_{\pm 4.37}$ | 682 000 $_{\pm 102\,274}$ |
| Length | 69.05 $_{\pm 4.15}$ | 63.84 $_{\pm 4.12}$ | 76.37 $_{\pm 4.46}$ | 672 667 $_{\pm 71\,760}$ |
| Word Rarity | 66.74 $_{\pm 3.48}$ | 62.48 $_{\pm 3.52}$ | 74.12 $_{\pm 4.12}$ | 664 666 $_{\pm 72\,394}$ |
| FRE | 68.05 $_{\pm 5.22}$ | 62.53 $_{\pm 4.98}$ | 75.32 $_{\pm 5.88}$ | 709 333 $_{\pm 105\,524}$ |
| Random | 68.07 $_{\pm 4.92}$ | 63.08 $_{\pm 4.95}$ | 75.21 $_{\pm 5.40}$ | 679 333 $_{\pm 105\,388}$ |

Table 3: Performance of BERT-tiny across baseline and CL strategies. Perplexity is reported for the full dataset and separately for the simple (SL) and everyday language (EL) subsets. Sequential label-based curriculum achieves best overall and SL perplexity. No competence-based strategy shows consistent improvement over baselines. Reported values are mean and standard deviations across 15 runs. ♪ denotes significant changes.

two attention heads, an intermediate feed-forward of size 512, a batch size of eight, and a learning rate of $10^{-4}$. All models are trained until convergence, with early stopping based on validation loss. All experiments are repeated over 15 random seeds to ensure statistical robustness.

**Evaluation** We evaluate model performance using overall perplexity as well as SL and EL subset perplexities. This helps us assess general improvements as well as register-specific gains. Our baselines include models trained with random sampling: one on everyday language only (Baseline EL), the other on a uniform mix (Baseline SL+EL).

## 4 Curriculum Learning Results

We summarise the final performance of the BERT-tiny model across all training strategies in Table 3, focusing on overall, SL, and EL perplexity, as loss values are less informative. We compare each strategy against a primary baseline (Baseline SL+EL), trained on SL+EL using random data sampling, with results averaged over 15 seeds. To assess the statistical significance of our results, we apply a one-sided Wilcoxon signed-rank test for symmetric distributions, and a one-sided median bootstrap test otherwise. All $p$-values are adjusted using the Holm-Bonferroni method within each experiment family (baseline, label-based CL, competence-based CL), using $\alpha = 0.05$ and directional hypotheses. Appendix C details the directional hypotheses and the corresponding adjusted $p$-values.

**Does merely adding simple language to the training data improve model performance?** The results provide a clear but mixed answer. Comparing Baseline SL+EL to Baseline EL, we see a significant improvement in EL perplexity but no improvement in neither overall nor SL perplexity.

**Can simple language effectively guide CL?** We find clear evidence in favour of simple language guiding CL – provided that the sampling strategy is right. Among the label-based CL strategies, only the sequential variant significantly improves overall as well as SL perplexity – achieving the best scores across all strategies. Incremental improves EL perplexity, but not overall performance. To show that the improvements of the sequential strategy are not accidental, we also test its anti strategy (i.e. starting training on EL, then progressing with SL): it performs similarly to Baseline EL and yields significantly worse EL perplexity than Baseline SL+EL. Both incremental and sequential strategies require more updates than Baseline SL+EL to reach these improvements.

**Are shallow text features sufficient to guide competence-based CL?** We have a negative answer to this question. Across all three competence-based difficulty measures, we observe no significant improvement in perplexity compared to Baseline SL+EL. The random strategy further suggests that neither simply increasing the dataset size nor imposing an order on shallow features leads to better model performance.

Figure 1: Distribution of sentence-level difficulty heuristics for SL and EL. None of the heuristics cleanly separates the two classes.

|     | SL      | EL      |
| --- | ------- | ------- |
| SL  | 100%    | 96.67%  |
| EL  | 86.06%  | 100%    |

Table 4: Vocabulary overlap between classes. Over 80% of EL's vocabulary is also present in SL, showing high lexical similarity.

## 5 Discussion

In this section we discuss the implications of the results from the previous section with regards to our three research questions.

**Learning across registers: asymmetries and interference** The surprisingly strong performance of Baseline EL on the SL subset suggests that EL may implicitly cover much of the SL distribution, possibly due to the compositionality of language. However, simply adding SL to the randomly ordered training data does not improve overall performance – and while it significantly improves EL perplexity, it worsens performance on SL itself. This asymmetry hints at a negative interference effect as observed in multilingual model training (Wang et al., 2020): though both classes stem from the same language, they might be different enough to cause gradient conflicts when used in the same dataset. These findings emphasise that learning patterns across language registers are not symmetric, and underscore the importance of evaluating perplexity for different subsets.

**Structure matters: the effectiveness of label-based curricula** Models only benefit from SL when introduced in a structured way. Sequential label-based curricula, where training begins with SL before using EL, consistently outperform other strategies in overall and SL perplexity. This aligns with the idea that simplified input can serve as a scaffold, supporting the acquisition of more complex patterns. While the effect mirrors principles observed in human learning, the underlying reason why structured exposure aids generalisation may differ in MLM.

**The limits of difficulty heuristics** Competence-based curricula using shallow difficulty heuristics show no clear advantage over random strategies. While this supports prior findings by Surkov et al. (2022), our analysis offers further insight. Figure 1 shows histograms comparing the distribution of shallow heuristics in SL and EL and Table 1 illustrates some examples. While it is plausible that EL has samples at the "easy" extremes, as not every sentence in everyday language is necessarily complex, we also observe SL examples at the "complex" extremes. Assuming that SL represents text that is easier to understand for humans, this highlights that the difficulty heuristics fail to meaningfully separate the two classes.

**Future Directions** We find that while shallow difficulty heuristics do not suffice to guide CL, the information encoded in the language classes does. Despite high lexical overlap and comparable size (Tables 2 and 4), simple language may offer more than surface-level simplicity. Prior work has shown that both humans and neural models benefit from regular, compositional input (Galke et al., 2024) and simple language might reflect just that through syntactic consistency or clearer discourse structure. Future work could explore how such compositional features manifest in simple language, and whether they can be modelled or annotated as difficulty signals – enabling broader and more effective CL strategies in MLM pre-training.

## 6 Conclusion

We examined whether human-curated simple language can guide CL in MLM pre-training. Our results show that label-based curricula outperform both random baselines and competence-based approaches relying on shallow difficulty heuristics. While the two language classes show high lexical and distributional overlap, their ordering – particularly when first training on simple language before moving to everyday language – leads to significant gains in model performance. This suggests that human intuition about linguistic difficulty provides more effective structure for CL than traditional surface-level heuristics.

# References

Alessia Battisti, Dominik Pfütze, Andreas Säuberli, Marek Kostrzewa, and Sarah Ebling. 2020. A Corpus for Automatic Readability Assessment and Text Simplification of German. In *Proceedings of the 12th Language Resources and Evaluation Conference, LREC 2020*, pages 3302–3311. European Language Resources Association.

Yoshua Bengio, Jérôme Louradour, Ronan Collobert, and Jason Weston. 2009. Curriculum Learning. In *Proceedings of the 26th International Conference On Machine Learning, ICML 2009*, pages 41–48.

William Coster and David Kauchak. 2011. Simple English Wikipedia: A New Text Simplification Task. In *Proceedings of the 49th Annual Meeting of the Association for Computational Linguistics: Human Language Technologies, ACL 2011*, volume 3, pages 665–669. Association for Computational Linguistics.

Jeffrey L. Elman. 1993. Learning and development in neural networks: the importance of starting small. *Cognition*, 48(1):71–99.

Rudolph Flesch. 1948. A New Readability Yardstick. *Journal of Applied Psychology*, 32:221–233.

Lukas Galke, Yoav Ram, and Limor Raviv. 2024. What Makes a Language Easy to Deep-Learn? Deep Neural Networks and Humans Similarly Benefit from Compositional Structure. *Nature Communications*, 15(1):10816.

Philip A. Huebner, Elior Sulem, Cynthia Fisher, and Dan Roth. 2021. BabyBERTa: Learning More Grammar With Small-Scale Child-Directed Language. In *Proceedings of the 25th Conference on Computational Natural Language Learning, CoNLL 2021*, pages 624–646. Association for Computational Linguistics.

David Kauchak. 2013. Improving Text Simplification Language Modeling Using Unsimplified Text Data. In *Proceedings of the 51st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers), ACL 2013*, volume 1, pages 1537–1546.

Cao Liu, Shizhu He, Kang Liu, and Jun Zhao. 2018. Curriculum Learning for Natural Answer Generation. In *Proceedings of the Twenty-Seventh International Joint Conference on Artificial Intelligence, IJCAI 2018*, pages 4223–4229. ijcai.org.

Evan Lucas, Dylan Gaines, Tagore Rao Kosireddy, Kevin Li, and Timothy C. Havens. 2024. Using Curriculum Masking Based on Child Language Development to Train a Large Language Model with Limited Training Data. In *The 2nd BabyLM Challenge at the 28th Conference on Computational Natural Language Learning*, pages 221–228. Association for Computational Linguistics.

Tasnim Mohiuddin, Philipp Koehn, Vishrav Chaudhary, James Cross, Shruti Bhosale, and Shafiq Joty. 2022. Data Selection Curriculum for Neural Machine Translation. In *Findings of the Association for Computational Linguistics: EMNLP 2022*, pages 1569–1582. Association for Computational Linguistics.

Aaron Mueller and Tal Linzen. 2023. How to Plant Trees in Language Models: Data and Architectural Effects on the Emergence of Syntactic Inductive Biases. In *Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers), ACL 2023*, pages 11237–11252. Association for Computational Linguistics.

Koichi Nagatsuka, Clifford Broni-Bediako, and Masayasu Atsumi. 2021. Pre-Training a BERT with Curriculum Learning by Increasing Block-Size of Input Text. In *Proceedings of the 12th International Conference on Recent Advances in Natural Language Processing, RANLP 2021*, pages 989–996. INCOMA Ltd., Shoumen, Bulgaria.

Emmanouil Antonios Platanios, Otilia Stretcu, Graham Neubig, Barnabas Poczos, and Tom M. Mitchell. 2019. Competence-Based Curriculum Learning for Neural Machine Translation. 1:1162–1172.

Leonardo Ranaldi, Giulia Pucci, and Fabio Massimo Zanzotto. 2023. Modeling Easiness for Training Transformers with Curriculum Learning. In *Proceedings of the 14th International Conference on Recent Advances in Natural Language Processing, RANLP 2023*, pages 937–948. INCOMA Ltd., Shoumen, Bulgaria.

Maxim Surkov, Vladislav Mosin, and Ivan Yamshchikov. 2022. Do Data-based Curricula Work? In *Proceedings of the Third Workshop on Insights from Negative Results in NLP*, pages 119–128. Association for Computational Linguistics.

Xin Wang, Yudong Chen, and Wenwu Zhu. 2021. A Survey on Curriculum Learning. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 44(9).

Zirui Wang, Zachary C. Lipton, and Yulia Tsvetkov. 2020. On Negative Interference in Multilingual Models: Findings and A Meta-Learning Treatment. In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing, EMNLP 2020*, pages 4438–4450. Association for Computational Linguistics.

Wojciech Zaremba and Ilya Sutskever. 2015. Learning to Execute. *Preprint*, arXiv:1410.4615.

## A  Implementation Details

We provide our implementation details for the competence-based CL strategy, where each training sample is assigned a difficulty score and the dataset is sorted accordingly. A predefined competence function then controls the fraction of data available at each training step $t$, gradually increasing the difficulty over time. Following Platanios et al. (2019), we adopt the square-root based competence function, which they found to be most effective:

$$c_{sqrt}(t) = \min(1, \sqrt{\frac{t(1 - c_0^2)}{T}}) \in [0, 1],$$

where $c_0$ denotes the initial competence at $t = 0$ and $T$ is the total number of steps in the CL phase. In our experiments, we observed that shorter competence phases tend to yield better results than longer ones. We pick $T = 50\,000$ and $c_0 = 0.05$ as function parameters. The size of the training dataset is updated every $5\,000$ steps depending on the current function value.

## B  Difficulty Heuristics

In our work we consider three popular heuristics to measure the difficulty of text for global, data-level curriculum learning (cf. Platanios et al. (2019) or Ranaldi et al. (2023)). Let $S$ be a sentence, represented by a finite sequence of words $(w_1, w_2, \ldots, w_m)$. The first heuristic, sentence length, is defined by the number of words in the sentence:

$$\text{length}(S) = |S|.$$

Next, we use the word rarity metric as proposed by Platanios et al. (2019), but normalise it by the number of words to remove its strong correlation with the sentence length:

$$\text{word rarity}(S) = -\frac{1}{|S|} \sum_{w \in S} \log\left(\frac{\text{count}_c(w)}{N}\right),$$

where $N$ denotes the size of the vocabulary of the corpus and $\text{count}_c(w)$ the number of times $w$ appeared in the corpus. Last, we present the Flesch Reading Ease (FRE) score as defined by Flesch (1948). It is designed to evaluate the readability of text and to return a score between 0 and 100:

$$\text{FRE}(S) = 206.835 - 1.015 \times \text{ASL} - 84.6 \times \text{ASW},$$

where ASL denotes the average sentence length, which is always the actual sentence length since we

| Strategy | PPL | SL PPL | EL PPL |
|---|---|---|---|
| Baseline SL+EL | .445 (w) | .996 (w) | **.004** (w) |
| Incremental | .598 (b) | 1.00 (w) | **.008** (w) |
| Sequential | **.019** (w) | **.001** (w) | .126 (w) |
| Anti-Sequential | .252 (b) | 1.00 (w) | **.008** (w) |
| Length | .890 (w) | .977 (w) | .899 (w) |
| Word Rarity | .890 (b) | .977 (w) | .718 (w) |
| FRE | .779 (w) | .977 (w) | .899 (w) |
| Random | .779 (w) | .977 (w) | .899 (w) |

Table 5: Adjusted $p$-values for all statistical tests for the models' performance on overall perplexity (PPL), simple language perplexity (SL PPL), and everyday language perplexity (EL PPL). We choose $\alpha = 0.05$ and boldface all significant results. We further indicate which one-sided test was run: (w) Wilcoxon signed-rank test or (b) boostrap median test.

only evaluate single sentences, and ASW denotes the average syllables per word. Since the FRE was designed to evaluate text samples of 100 words, we can encounter negative FRE scores which are outside the originally defined range.

## C  Details on the Significance Tests

Table 5 reports the adjusted $p$-values for all strategies, assessing their performance relative to relevant baselines. For each comparison, we applied a one-sided test based on our directional hypotheses: (1) whether adding SL (Baseline SL+EL) *improves* over the baseline trained with EL (Baseline EL); (2) whether label-based curricula (Incremental and Sequential) *improve* over the full baseline (Baseline SL+EL); (3) whether Anti-Sequential *hurts* performance compared to Baseline SL+EL; and (4) whether competence-based strategies (Length, Word Rarity, FRE, Random) *improve* over the Baseline SL+EL.

# FActBench: A Benchmark for Fine-grained Automatic Evaluation of LLM-Generated Text in the Medical Domain

**Anum Afzal**
Technical University of Munich
anum.afzal@tum.de

**Juraj Vladika**
Technical University of Munich
juraj.vladika@tum.de

**Florian Matthes**
Technical University of Munich
matthes@tum.de

## Abstract

Large Language Models tend to struggle when dealing with specialized domains. While all aspects of evaluation hold importance, factuality is the most critical one. Similarly, reliable fact-checking tools and data sources are essential for hallucination mitigation. We address these issues by providing a comprehensive Fact-checking Benchmark FActBench covering four generation tasks and six state-of-the-art Large Language Models (LLMs) for the Medical domain. We use two state-of-the-art Fact-checking techniques: Chain-of-Thought (CoT) Prompting and Natural Language Inference (NLI). Our experiments show that the fact-checking scores acquired through the Unanimous Voting of both techniques correlate best with Domain Expert Evaluation.

## 1 Introduction

In the quickly evolving era of Natural Language Processing (NLP), Large Language Models (LLMs) are making their way into almost all use cases and domains. In most tasks, they have shown tremendous generative capabilities and a good understanding of text. However, they still tend to hallucinate in critical domains like the Medical domain. Contemporary LLMs are typically evaluated against general benchmarks and their assessment of the Medical domain is usually lacking. While it is essential to mitigate hallucinations, as a first step some reliable automatic fact-checking indicators are needed (Clusmann et al., 2023). The field of automatic fact-checking in LLMs is rapidly developing making it essential to find trustworthy techniques and data sources.

The state-of-the-art techniques for Automatic Fact Checking include Natural Language Inference (NLI) (Mor-Lan and Levi, 2024; Akhtar et al., 2024) using DeBERTa (He et al., 2021), or through Chain-of-thought (CoT) (Wei et al., 2022) by using an LLM as a judge (Zheng et al., 2023). Given the importance of Factual correctness in a critical domain such as medicine, it is helpful to rely on more than one technique for Fact-checking. Therefore, we explore the idea of Unanimous Voting such that an atomic fact is only considered to be factually correct if it is supported by both techniques.

Hallucinations can generally be divided into input-conflicting, context-conflicting, and fact-conflicting (Zhang et al., 2023). The focus of our work lies in fact-conflicting, which is hallucination, where facts in output contradict the world knowledge. Additionally, our work builds on top of FActScore, a CoT-based approach for fact-checking. We adapt it to support user-provided grounding documents, making it suitable for tasks like RAG and Summarization. We present an automatic Fact-Checking Benchmark **FActBench**[1] with the following contributions:

- We fact-check six contemporary LLMs using Atomic Facts (Min et al., 2023) on four generations tasks: Text Summarization, Lay Summarization, Retrieval Augmented Generation (RAG), and Open-ended Generation.

- We compare Intrinsic (Grounding Document) and Extrinsic (Wikipedia Dump) Fact-checking techniques in our experiments.

- We evaluate NLI, CoT as well as Unanimous Voting (UnVot) for the final prediction using domain expert evaluations as reference.

Details about all the datasets we use can be found in their original papers, including appropriate licenses and terms of use.

## 2 Related Work

Hallucinations are a common problem in Natural Language Generation (NLG) tasks such as abstrac-

---

[1]Code for FActBench can be found at github.com/jvladika/FactSumm/

tive text summarization, generative question answering, or dialogue generation (Ji et al., 2023). Detecting hallucinations is tied to the problem of measuring the factuality of model output (Augenstein et al., 2023; Zhao et al., 2024). Hallucinations can be detected with approaches looking at the uncertainty in models' logits (Varshney et al., 2023) or with approaches that fact-check model output over external knowledge sources (Chern et al., 2023).

Some recent works approached evaluation with question answering (Scialom et al., 2021) or NLI (Utama et al., 2022). Most recent methods leverage LLMs by querying them with prompts that directly ask for a score, like G-Eval (Liu et al., 2023a), or evaluate the generated text with the veracity of its atomic facts, like FActScore (Min et al., 2023). Fadeeva et al. (2024) develop a method that does not require external knowledge for fact-checking as they leverage token-level uncertainty to identify the potentially factually incorrect generated section in the output. Similarly, Sankararaman et al. (2024) introduces Provenance, a technique that uses NLI models to check if the RAG output is factually correct with reference to context. Lastly, Chen et al. (2024) present FactCHD, a benchmarking for fact-conflicting hallucination detection for the General, Scientific, Health, and COVID-19 domains.

## 3 FActBench: Benchmark

In our Benchmark, we use two SotA techniques, NLI and CoT, to evaluate 6 models on 4 different tasks. We follow the approach introduced by Min et al. (2023) to break all generations into a list of atomic facts which are then used for fact-checking. Since all our tasks with the exception of Open Generation, use a source document for grounding, we opt for a hybrid approach such that we first perform fact-checking using an intrinsic approach, followed by an extrinsic one.[2] The latter only performs evaluation on atomic facts that have been marked as hallucinations in the first step. We employ such an approach because it is possible for an atomic fact to be factually correct as per the world knowledge, even if it is not supported by the grounding document. We show this methodology in Figure 1 that illustrates how different fact-checking techniques and data sources interact with each other.

### 3.1 Techniques

**Baseline: FActScore** As a baseline, we first report on task performance using the established FActScore metric, following their external checks on Wikipedia with no grounding document. The reason we use it is its popularity in papers involving generative NLP tasks in the last couple of years (Dhuliawala et al., 2024; Chang et al., 2024; Huang et al., 2025). Later, our goal is to show that the combination of methods we use instead of raw FActScore lead to a more faithful evaluation framework and a better alignment with human scores.

**Natural Language Inference (NLI):** We utilize NLI as the first evaluation method. NLI aims to predict the logical relation between a premise and a hypothesis, including entailment, contradiction, and a neutral stance. We use the generated answer as the premise and the reference answer as the hypothesis. The intuition behind this approach is that a good answer should logically entail the reference. NLI has been applied for evaluating the quality of summaries and text generation (Mishra et al., 2021; Laban et al., 2022; Steen et al., 2023).

Following this approach, we use DeBERTa-v3 (He et al., 2023), shown to work well with NLI and reasoning tasks. We use the version *Tasksource*, fine-tuned on a wide array of NLI & classification datasets, which works well with long inputs (Sileo, 2023).[3] We take *entailment* predictions as a sign of the atomic fact being supported by the original text and *contradiction* as a sign of hallucination. We additionally check the contradicting atomic facts in an extrinsic way, by predicting their NLI class with the relevant Wikipedia context as the hypothesis.

**Chain-of Thought (CoT) Prompting:** For evaluation using Chain-of-Thought Prompting, we adapted FActScore, an existing CoT-based fact-checking tool. This technique is suitable for open-ended generation and uses a Wikipedia dump as the knowledge source. FActScore supports extrinsic fact-checking by retrieving the most relevant passages from Wikipedia using user-defined topics. We adapt FActScore to support external documents as the basis for fact-checking. This "topic" should be the name of a real Wikipedia article, from which the relevant passages are retrieved. We also include a LLM-based topic generator so it is not required to manually define the topic when evaluating using

---

[2]In factuality evaluation, *intrinsic* hallucinations are those that contradict the reference document, while *extrinsic* hallucinations are those that contradict the external world knowledge.

[3]https://huggingface.co/tasksource/deberta-base-long-nli

Figure 1: Block Diagram depicting how different fact-checking techniques interact with different data sources. Chain-of-Thought uses an LLM whereas Natural Language Inference uses a small LM as the backbone.

passages from the Wikipedia dump. We use GPT-4o mini as the backbone of FActScore+, which serves as a compromise between cost and quality.

**Unanimous Voting (UnVot):** To produce a reliable fact-checking approach, we explore the idea of Unanimous Voting. This means we only consider an atomic fact to be correct if both NLI and CoT support it. This technique is especially useful for applications where high precision is needed.

**Human Evaluation:** We evaluate CoT, NLI, and UnVot techniques by correlating to domain expert judgment. We recruited 8 in-house employed individuals with a medical background to serve as annotators. A random subset of 80 generations (20 per task) was manually annotated such that each generation was evaluated by two annotators. They were instructed to follow the same hybrid, using both the original article and Wikipedia as a basis for fact-checking. Annotators were asked to assign a score between 1 and 100 to the generation estimating the factual correctness of the text.

### 3.2 Tasks

We include four tasks in our Benchmark, including Text Summarization, Lay Summarization, Retrieval Augmented Generation, and Open-ended Generation. The prompts used for all four tasks are shown in Appendix A. We summarize the datasets used for the tasks in Table 1 and discuss them below. All the datasets can be found in their respective original papers, together with appropriate licenses.

**Text Summarization.** This task refers to the ability of an LLM to summarize a long scientific article into a summary. We used 1000 random samples from the PubMed Summarization dataset (Cohan et al., 2018), which is derived from the

original PubMed dump.

**Lay Summarization.** Contrary to normal text summarization, Lay Summarization refers to the model's ability to create a layman summary of biomedical articles. We use 1000 random samples from the PLOS dataset introduced by Goldsack et al. (2022).

**Retrieval Augmented Generation (RAG).** We use BioASQ-QA (Krithara et al., 2023), a biomedical question answering (QA) dataset designed to reflect the real information needs of biomedical experts. The questions are written by experts and evidence comes from PubMed. We use the *summary* subset – 1130 questions paired with human-selected evidence snippets from PubMed and human-written "ideal answers" based on those snippets. We use the gold snippets as input to an LLM and prompt it to generate an answer to the given question, thus simulating a **RAG** pipeline.

**Open-ended Generation.** In this setting, no context is used and the model is prompted to generate an answer based on its knowledge. We again use the BioASQ dataset from the RAG task – we take the 1130 questions and use them as input to an LLM by prompting it to answer the question.

| Task | Dataset | #Source W | #Gen W |
|------|---------|-----------|--------|
| Summ | PubMed | 3,053.9 | 256 |
| Lay Summ | PLOS | 6,696.8 | 256 |
| RAG | BioASQ-QA | 351.9 | 116.5 |
| Gen | BioASQ-QA | 351.9 | default |

Table 1: Average word count of articles (#W) and # generation tokens (#Gen W) during inference for tasks with respective datasets. Summ = Text Summarization, Lay Summ = Layman Summarization, RAG = Retrieval Augmented Generation, Gen = Open-ended Generation.

| Models | Summarization | | | Lay Summarization | | | RAG (QA) | | | Open-ended Gen | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | CoT | NLI | UnVot | CoT | NLI | UnVot | CoT | NLI | UnVot | CoT | NLI | UnVot |
| *FActBench (Grounding Document)* | | | | | | | | | | | | |
| GPT-4o mini | 95.8 | 77.4 | 86.6 | 95.4 | 94.8 | 95.1 | 25.4 | **77.7** | 51.5 | 44.5 | **50.4** | 47.5 |
| Llama3.1 8b | 95.3 | 87.8 | 85.28 | 95.4 | 93.5 | 94.4 | 35.6 | 76.7 | 56.1 | 74.4 | 35.6 | 55.0 |
| Llama3.1 70b | **96.52** | 84.59 | 82.84 | 96.1 | 94.1 | 95.1 | 31.3 | 76.3 | 53.8 | 37.3 | 46.5 | 41.8 |
| Mistral 7b | 95.8 | 82.55 | 80.38 | 96.3 | 97.32 | 94.75 | 82.9 | 73.1 | 78.0 | 80.9 | 32.5 | 56.6 |
| Mixtral 8 x 7b | 95.2 | 87.86 | **95.5** | 96.5 | 97.0 | 95.0 | **88.2** | 75.0 | **81.6** | **85.4** | 36.9 | **61.1** |
| Gemma 9b | 84.55 | 71.95 | 68.77 | 82.94 | 80.65 | 75.48 | 35.8 | 44.0 | 43.7 | 54.1 | 30.5 | 28.0 |
| *FActBench (Grounding Document + Wikipedia)* | | | | | | | | | | | | |
| GPT-4o mini | 96.8 | 82.6 | 80.4 | 96.2 | 96.6 | 93.4 | 97.3 | **78.2** | **76.4** | **95.8** | **51.4** | **50.3** |
| Llama3.1 8b | 96.4 | 88.85 | 86.25 | 96.5 | 94.2 | 91.5 | 98.2 | 77.1 | 76.1 | 79.3 | 36.7 | 32.1 |
| Llama3.1 70b | 97.27 | 85.71 | 83.9 | 97.0 | 94.8 | 92.0 | 97.2 | 76.8 | 75.1 | 90.9 | 47.7 | 45.9 |
| Mistral 7b | 96.51 | 83.59 | 81.34 | **97.83** | **96.7** | **94.93** | **98.6** | 73.5 | 72.7 | 92.1 | 33.2 | 31.9 |
| Mixtral 8 x 7b | 96.9 | 88.68 | **86.24** | 97.5 | 97.2 | 95.1 | 97.7 | 75.3 | 74.0 | 93.0 | 37.8 | 36.5 |
| Gemma 9b | 93.03 | 74.46 | 70.99 | 91.11 | 81.68 | 76.43 | 97.4 | 45.0 | 44.6 | 80.1 | 31.5 | 28.8 |
| *Baseline: FActScore (Wikipedia)* | | | | | | | | | | | | |
| GPT-4o mini | | 51.34 | | | 52.6 | | | 19.4 | | | 41.4 | |
| Llama3.1 8b | | 43.97 | | | 49.4 | | | 25.3 | | | 71.3 | |
| Llama3.1 70b | | 50.08 | | | 48.8 | | | 24.0 | | | 34.8 | |
| Mistral 7b | | 46.11 | | | 50.02 | | | 61.1 | | | 78.4 | |
| Mixtral 8 x 7b | | 49.71 | | | 51.00 | | | **64.5** | | | **81.6** | |
| Gemma 9b | | **53.54** | | | **54.56** | | | 44.0 | | | 52.0 | |

Table 2: Factchecking scores of six LLMs on four tasks using Chain-of-Thought (CoT) prompting, Natural Language Inference (NLI), and Unanimous voting (UnVot). We show scores by incorporating two different knowledge sources.

## 3.3 Models

We include six LLMs in our experiments including Llama3.1 8b (Dubey et al., 2024) Llama3.1 70b, Mistral 7b (Jiang et al., 2023), Mixtral 8x7b (Jiang et al., 2024a), Gemma2 9b (Team et al., 2024) and lastly, closed-source GPT-4o mini. We provide the checkpoints and technical details in Appendix B.

## 4 Results & Discussion

## 4.1 Correlation with Human Evaluation

Before discussing the benchmark results, we check the effectiveness of the techniques used. We performed human evaluation using the process described in subsection 3.1. The average fact-checking scores using the baseline, 3 techniques, as well domain expert annotations are in Table 3. The final Cohen's inter-annotator agreement $\kappa$ is 0.75, which signifies substantial agreement. The baseline technique (FActScore) that uses only Wikipedia as the knowledge source severely underestimates the correctness of the generated text whereas the Chain-of-Thought technique that uses Grounding Document and Wikipedia overestimates it. Overall, it can be seen that our UnVot score derived through joint decisions of CoT and NLI correlates best with domain expert judgment. Still, it is important to point out that this holds true for the summarization, lay summarization, and RAG tasks, while the pure generation task best correlated with baseline FActScore system.

The high correlation of UnVot with human judgment is an important finding. Hiring human annotators, especially domain experts, can be a very expensive and time-consuming process. Having a metric that highly correlates with human scoring intuition can provide a good enough substitute for situations where finding human annotators is infeasible or impossible for certain labs, groups, and application use cases. A lot of focus of recent LLM research is put on aligning LLMs with human values and intuition (Wang et al., 2023), and recent LLM-as-judge evaluation metrics like G-Eval (Liu et al., 2023b), Prometheus (Kim et al., 2024), and TIGERScore (Jiang et al., 2024b) put a high emphasis on the correlation of their metrics with humans. As future work, it would be interesting to compare these metrics with UnVot as well, which we currently skip due to resource constraints.

## 4.2 Task and LLM Performance

We summarize the Fact-checking scores in Table 2, which show that the grounding helps LLMs to be more truthful. In terms of tasks, LLMs tend to hallucinate more when prompted to do open-ended

| Task | Baseline | CoT* | NLI* | UnVot* | Human |
|---|---|---|---|---|---|
| Summ | 54.81 | 96.87 | 85.41 | 83.45 | 84.0 |
| LaySumm | 52.5 | 97.6 | 91.09 | 88.94 | 88.7 |
| RAG | 38.43 | 100.0 | 83.04 | 83.04 | 87.3 |
| PureGen | 71.26 | 88.17 | 31.61 | 31.31 | 62.7 |

Table 3: Fact-checking scores on FActScore (Baseline), Chain-of-Thought (CoT), Natural Language Inference (NLI), Unanimous Voting (UnVot), and Domain Expert Evaluation (Human). * refers to final scores with intrinsic followed by extrinsic fact-checking.

generation in the medical domain. However, the performance on other grounding-based task show that given the correct context and supporting document, LLMs are good at understanding a complex domain such as the medical domain. Within each task, LLM performance is mostly uniform. As expected, Open-ended generation is the most challenging task, which is expected due to the LLM using its internal knowledge to answer questions, which can lead to hallucinations. Lay summarization was the most factually correct task, likely owing to the nature of lay text where simpler terms and phrasing is used, which reduces the possibility of mixing up complex scientific terms with one another, which would lead to hallucinations.

Surprisingly, we see no big difference in models with respect to their sizes. However, both `Mistral` and `Mixtral` lead the board for two summarization tasks. While `Mixtral` performs best for two QA tasks with only the grounding document, `GPT` comes on top after extrinsic checks, showing its high awareness of Wikipedia in pre-trained knowledge. Two `Llama` models come close to `Mixtral`, while `Gemma` performs the worst on all tasks.

## 5 Conclusion and Future Work

We present a Benchmark providing insights over contemporary LLMs across 4 tasks in the medical domain. We discuss Chain-of-Thought, Natural Language Inference, and Unanimous Voting as fact-checking techniques. Through Domain Expert Evaluation, we show the Unanimous Voting technique to be most reliable. We also explored the effectiveness of two knowledge sources, namely a Grounding Document and Wikipedia, for evaluation and found that using more than one knowledge source leads to an increase in factuality scores. Lastly, we found that LLMs are mostly factually incorrect for Open-ended generation in the medical domain and tend to be more faithful for tasks like

Summarization and RAG, where some context is provided to the LLM for generation. We envision our evaluation benchmark to be easily applied for fact-checking across other domains in future.

## Limitation

Due to the high computation costs, we use only one model as the backbone for each factuality evaluation technique. Even though we evaluated six Large Language Models on four diverse tasks, these tasks may not be enough to capture the entirety of LLM performance and the quickly evolving landscape of new models.

Additionally, our two evaluation techniques with NLI and FactScore+ CoT are not perfect and it is possible there were incorrect predictions of which facts were supported or refuted by evidence. Even though our manual inspection and human evaluation showed a good correlation with automated metrics, there will always be some mishaps and incorrect verdicts.

Finally, our approach relies on making numerous calls to the external API and to the Wikipedia dump database instance in case of extrinsic fact-checking, which can all slow down the overall pipeline. An alternative would have been running locally hosted open-source models, but this was out of our budget due to computational costs. Future work could explore these solutions and make the process faster.

## Ethics Statement

Throughout our experiments, we strictly adhere to the ACL Code of Ethics. The manual evaluation was performed by in-house annotators who received a full salary, and their annotation were stored anonymously, mitigating any privacy concerns. They were informed about the task and usability of data in the research. The goal of the research is to evaluate existing techniques and introduce a new technique that can be used for fact-checking LLM generated text on four tasks in the medical domain. We use the LLMs through inference using open-source dataset and do not include in any information in model weights. The discussions and results in this paper are meant to further promote research in the area of LLM Fact-checking as well as create more awareness about their applications in the medical domain. All scripts will be made available to the research community.

# References

Mubashara Akhtar, Michael Schlichtkrull, and Andreas Vlachos. 2024. Ev2r: Evaluating evidence retrieval in automated fact-checking. *Preprint*, arXiv:2411.05375.

Isabelle Augenstein, Timothy Baldwin, Meeyoung Cha, Tanmoy Chakraborty, Giovanni Luca Ciampaglia, David Corney, Renee DiResta, Emilio Ferrara, Scott Hale, Alon Halevy, Eduard Hovy, Heng Ji, Filippo Menczer, Ruben Miguez, Preslav Nakov, Dietram Scheufele, Shivam Sharma, and Giovanni Zagni. 2023. Factuality challenges in the era of large language models. *Preprint*, arXiv:2310.05189.

Yupeng Chang, Xu Wang, Jindong Wang, Yuan Wu, Linyi Yang, Kaijie Zhu, Hao Chen, Xiaoyuan Yi, Cunxiang Wang, Yidong Wang, et al. 2024. A survey on evaluation of large language models. *ACM transactions on intelligent systems and technology*, 15(3):1–45.

Xiang Chen, Duanzheng Song, Honghao Gui, Chenxi Wang, Ningyu Zhang, Yong Jiang, Fei Huang, Chengfei Lv, Dan Zhang, and Huajun Chen. 2024. Factchd: Benchmarking fact-conflicting hallucination detection. *Preprint*, arXiv:2310.12086.

I-Chun Chern, Steffi Chern, Shiqi Chen, Weizhe Yuan, Kehua Feng, Chunting Zhou, Junxian He, Graham Neubig, Pengfei Liu, et al. 2023. Factool: Factuality detection in generative ai–a tool augmented framework for multi-task and multi-domain scenarios. *arXiv preprint arXiv:2307.13528*.

Jan Clusmann, Fiona R Kolbinger, Hannah Sophie Muti, Zunamys I Carrero, Jan-Niklas Eckardt, Narmin Ghaffari Laleh, Chiara Maria Lavinia Löffler, Sophie-Caroline Schwarzkopf, Michaela Unger, Gregory P Veldhuizen, et al. 2023. The future landscape of large language models in medicine. *Communications medicine*, 3(1):141.

Arman Cohan, Franck Dernoncourt, Doo Soon Kim, Trung Bui, Seokhwan Kim, Walter Chang, and Nazli Goharian. 2018. A discourse-aware attention model for abstractive summarization of long documents. In *Proceedings of the 2018 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 2 (Short Papers)*, pages 615–621, New Orleans, Louisiana. Association for Computational Linguistics.

Shehzaad Dhuliawala, Mojtaba Komeili, Jing Xu, Roberta Raileanu, Xian Li, Asli Celikyilmaz, and Jason Weston. 2024. Chain-of-verification reduces hallucination in large language models. In *Findings of the Association for Computational Linguistics: ACL 2024*, pages 3563–3578, Bangkok, Thailand. Association for Computational Linguistics.

Abhimanyu Dubey, Abhinav Jauhri, Abhinav Pandey, Abhishek Kadian, Ahmad Al-Dahle, Aiesha Letman, Akhil Mathur, Alan Schelten, Amy Yang, Angela Fan, Anirudh Goyal, Anthony Hartshorn, Aobo Yang, Archi Mitra, Archie Sravankumar, Artem Korenev, Arthur Hinsvark, Arun Rao, Aston Zhang, Aurelien Rodriguez, Austen Gregerson, Ava Spataru, Baptiste Roziere, Bethany Biron, Binh Tang, Bobbie Chern, Charlotte Caucheteux, Chaya Nayak, Chloe Bi, Chris Marra, Chris McConnell, Christian Keller, Christophe Touret, Chunyang Wu, Corinne Wong, Cristian Canton Ferrer, Cyrus Nikolaidis, Damien Allonsius, Daniel Song, Danielle Pintz, Danny Livshits, David Esiobu, Dhruv Choudhary, Dhruv Mahajan, Diego Garcia-Olano, Diego Perino, Dieuwke Hupkes, Egor Lakomkin, Ehab AlBadawy, Elina Lobanova, Emily Dinan, Eric Michael Smith, Filip Radenovic, Frank Zhang, Gabriel Synnaeve, Gabrielle Lee, Georgia Lewis Anderson, Graeme Nail, Gregoire Mialon, Guan Pang, Guillem Cucurell, Hailey Nguyen, Hannah Korevaar, Hu Xu, Hugo Touvron, Iliyan Zarov, Imanol Arrieta Ibarra, Isabel Kloumann, Ishan Misra, Ivan Evtimov, Jade Copet, Jaewon Lee, Jan Geffert, Jana Vranes, Jason Park, Jay Mahadeokar, Jeet Shah, Jelmer van der Linde, Jennifer Billock, Jenny Hong, Jenya Lee, Jeremy Fu, Jianfeng Chi, Jianyu Huang, Jiawen Liu, Jie Wang, Jiecao Yu, Joanna Bitton, Joe Spisak, Jongsoo Park, Joseph Rocca, Joshua Johnstun, Joshua Saxe, Junteng Jia, Kalyan Vasuden Alwala, Kartikeya Upasani, Kate Plawiak, Ke Li, Kenneth Heafield, Kevin Stone, Khalid El-Arini, Krithika Iyer, Kshitiz Malik, Kuenley Chiu, Kunal Bhalla, Lauren Rantala-Yeary, Laurens van der Maaten, Lawrence Chen, Liang Tan, Liz Jenkins, Louis Martin, Lovish Madaan, Lubo Malo, Lukas Blecher, Lukas Landzaat, Luke de Oliveira, Madeline Muzzi, Mahesh Pasupuleti, Mannat Singh, Manohar Paluri, Marcin Kardas, Mathew Oldham, Mathieu Rita, Maya Pavlova, Melanie Kambadur, Mike Lewis, Min Si, Mitesh Kumar Singh, Mona Hassan, Naman Goyal, Narjes Torabi, Nikolay Bashlykov, Nikolay Bogoychev, Niladri Chatterji, Olivier Duchenne, Onur Çelebi, Patrick Alrassy, Pengchuan Zhang, Pengwei Li, Petar Vasic, Peter Weng, Prajjwal Bhargava, Pratik Dubal, Praveen Krishnan, Punit Singh Koura, Puxin Xu, Qing He, Qingxiao Dong, Ragavan Srinivasan, Raj Ganapathy, Ramon Calderer, Ricardo Silveira Cabral, Robert Stojnic, Roberta Raileanu, Rohit Girdhar, Rohit Patel, Romain Sauvestre, Ronnie Polidoro, Roshan Sumbaly, Ross Taylor, Ruan Silva, Rui Hou, Rui Wang, Saghar Hosseini, Sahana Chennabasappa, Sanjay Singh, Sean Bell, Seohyun Sonia Kim, Sergey Edunov, Shaoliang Nie, Sharan Narang, Sharath Raparthy, Sheng Shen, Shengye Wan, Shruti Bhosale, Shun Zhang, Simon Vandenhende, Soumya Batra, Spencer Whitman, Sten Sootla, Stephane Collot, Suchin Gururangan, Sydney Borodinsky, Tamar Herman, Tara Fowler, Tarek Sheasha, Thomas Georgiou, Thomas Scialom, Tobias Speckbacher, Todor Mihaylov, Tong Xiao, Ujjwal Karn, Vedanuj Goswami, Vibhor Gupta, Vignesh Ramanathan, Viktor Kerkez, Vincent Gonguet, Virginie Do, Vish Vogeti, Vladan Petrovic, Weiwei Chu, Wenhan Xiong, Wenyin Fu, Whitney Meers, Xavier Martinet, Xiaodong Wang, Xiaoqing Ellen Tan, Xinfeng Xie, Xuchao Jia, Xuewei Wang, Yaelle Goldschlag, Yashesh Gaur, Yasmine

Babaei, Yi Wen, Yiwen Song, Yuchen Zhang, Yue Li, Yuning Mao, Zacharie Delpierre Coudert, Zheng Yan, Zhengxing Chen, Zoe Papakipos, Aaditya Singh, Aaron Grattafiori, Abha Jain, Adam Kelsey, Adam Shajnfeld, Adithya Gangidi, Adolfo Victoria, Ahuva Goldstand, Ajay Menon, Ajay Sharma, Alex Boesenberg, Alex Vaughan, Alexei Baevski, Allie Feinstein, Amanda Kallet, Amit Sangani, Anam Yunus, Andrei Lupu, Andres Alvarado, Andrew Caples, Andrew Gu, Andrew Ho, Andrew Poulton, Andrew Ryan, Ankit Ramchandani, Annie Franco, Aparajita Saraf, Arkabandhu Chowdhury, Ashley Gabriel, Ashwin Bharambe, Assaf Eisenman, Azadeh Yazdan, Beau James, Ben Maurer, Benjamin Leonhardi, Bernie Huang, Beth Loyd, Beto De Paola, Bhargavi Paranjape, Bing Liu, Bo Wu, Boyu Ni, Braden Hancock, Bram Wasti, Brandon Spence, Brani Stojkovic, Brian Gamido, Britt Montalvo, Carl Parker, Carly Burton, Catalina Mejia, Changhan Wang, Changkyu Kim, Chao Zhou, Chester Hu, Ching-Hsiang Chu, Chris Cai, Chris Tindal, Christoph Feichtenhofer, Damon Civin, Dana Beaty, Daniel Kreymer, Daniel Li, Danny Wyatt, David Adkins, David Xu, Davide Testuggine, Delia David, Devi Parikh, Diana Liskovich, Didem Foss, Dingkang Wang, Duc Le, Dustin Holland, Edward Dowling, Eissa Jamil, Elaine Montgomery, Eleonora Presani, Emily Hahn, Emily Wood, Erik Brinkman, Esteban Arcaute, Evan Dunbar, Evan Smothers, Fei Sun, Felix Kreuk, Feng Tian, Firat Ozgenel, Francesco Caggioni, Francisco Guzmán, Frank Kanayet, Frank Seide, Gabriela Medina Florez, Gabriella Schwarz, Gada Badeer, Georgia Swee, Gil Halpern, Govind Thattai, Grant Herman, Grigory Sizov, Guangyi, Zhang, Guna Lakshminarayanan, Hamid Shojanazeri, Han Zou, Hannah Wang, Hanwen Zha, Haroun Habeeb, Harrison Rudolph, Helen Suk, Henry Aspegren, Hunter Goldman, Ibrahim Damlaj, Igor Molybog, Igor Tufanov, Irina-Elena Veliche, Itai Gat, Jake Weissman, James Geboski, James Kohli, Japhet Asher, Jean-Baptiste Gaya, Jeff Marcus, Jeff Tang, Jennifer Chan, Jenny Zhen, Jeremy Reizenstein, Jeremy Teboul, Jessica Zhong, Jian Jin, Jingyi Yang, Joe Cummings, Jon Carvill, Jon Shepard, Jonathan McPhie, Jonathan Torres, Josh Ginsburg, Junjie Wang, Kai Wu, Kam Hou U, Karan Saxena, Karthik Prasad, Kartikay Khandelwal, Katayoun Zand, Kathy Matosich, Kaushik Veeraraghavan, Kelly Michelena, Keqian Li, Kun Huang, Kunal Chawla, Kushal Lakhotia, Kyle Huang, Lailin Chen, Lakshya Garg, Lavender A, Leandro Silva, Lee Bell, Lei Zhang, Liangpeng Guo, Licheng Yu, Liron Moshkovich, Luca Wehrstedt, Madian Khabsa, Manav Avalani, Manish Bhatt, Maria Tsimpoukelli, Martynas Mankus, Matan Hasson, Matthew Lennie, Matthias Reso, Maxim Groshev, Maxim Naumov, Maya Lathi, Meghan Keneally, Michael L. Seltzer, Michal Valko, Michelle Restrepo, Mihir Patel, Mik Vyatskov, Mikayel Samvelyan, Mike Clark, Mike Macey, Mike Wang, Miquel Jubert Hermoso, Mo Metanat, Mohammad Rastegari, Munish Bansal, Nandhini Santhanam, Natascha Parks, Natasha White, Navyata Bawa, Nayan Singhal, Nick Egebo, Nicolas Usunier, Nikolay Pavlovich Laptev, Ning Dong, Ning Zhang, Norman Cheng, Oleg

Chernoguz, Olivia Hart, Omkar Salpekar, Ozlem Kalinli, Parkin Kent, Parth Parekh, Paul Saab, Pavan Balaji, Pedro Rittner, Philip Bontrager, Pierre Roux, Piotr Dollar, Polina Zvyagina, Prashant Ratanchandani, Pritish Yuvraj, Qian Liang, Rachad Alao, Rachel Rodriguez, Rafi Ayub, Raghotham Murthy, Raghu Nayani, Rahul Mitra, Raymond Li, Rebekkah Hogan, Robin Battey, Rocky Wang, Rohan Maheswari, Russ Howes, Ruty Rinott, Sai Jayesh Bondu, Samyak Datta, Sara Chugh, Sara Hunt, Sargun Dhillon, Sasha Sidorov, Satadru Pan, Saurabh Verma, Seiji Yamamoto, Sharadh Ramaswamy, Shaun Lindsay, Shaun Lindsay, Sheng Feng, Shenghao Lin, Shengxin Cindy Zha, Shiva Shankar, Shuqiang Zhang, Shuqiang Zhang, Sinong Wang, Sneha Agarwal, Soji Sajuyigbe, Soumith Chintala, Stephanie Max, Stephen Chen, Steve Kehoe, Steve Satterfield, Sudarshan Govindaprasad, Sumit Gupta, Sungmin Cho, Sunny Virk, Suraj Subramanian, Sy Choudhury, Sydney Goldman, Tal Remez, Tamar Glaser, Tamara Best, Thilo Kohler, Thomas Robinson, Tianhe Li, Tianjun Zhang, Tim Matthews, Timothy Chou, Tzook Shaked, Varun Vontimitta, Victoria Ajayi, Victoria Montanez, Vijai Mohan, Vinay Satish Kumar, Vishal Mangla, Vítor Albiero, Vlad Ionescu, Vlad Poenaru, Vlad Tiberiu Mihailescu, Vladimir Ivanov, Wei Li, Wenchen Wang, Wenwen Jiang, Wes Bouaziz, Will Constable, Xiaocheng Tang, Xiaofang Wang, Xiaojian Wu, Xiaolan Wang, Xide Xia, Xilun Wu, Xinbo Gao, Yanjun Chen, Ye Hu, Ye Jia, Ye Qi, Yenda Li, Yilin Zhang, Ying Zhang, Yossi Adi, Youngjin Nam, Yu, Wang, Yuchen Hao, Yundi Qian, Yuzi He, Zach Rait, Zachary DeVito, Zef Rosnbrick, Zhaoduo Wen, Zhenyu Yang, and Zhiwei Zhao. 2024. The llama 3 herd of models. *Preprint*, arXiv:2407.21783.

Ekaterina Fadeeva, Aleksandr Rubashevskii, Artem Shelmanov, Sergey Petrakov, Haonan Li, Hamdy Mubarak, Evgenii Tsymbalov, Gleb Kuzmin, Alexander Panchenko, Timothy Baldwin, Preslav Nakov, and Maxim Panov. 2024. Fact-checking the output of large language models via token-level uncertainty quantification. In *Findings of the Association for Computational Linguistics: ACL 2024*, pages 9367–9385, Bangkok, Thailand. Association for Computational Linguistics.

Tomas Goldsack, Zhihao Zhang, Chenghua Lin, and Carolina Scarton. 2022. Making science simple: Corpora for the lay summarisation of scientific literature. In *Proceedings of the 2022 Conference on Empirical Methods in Natural Language Processing*, pages 10589–10604, Abu Dhabi, United Arab Emirates. Association for Computational Linguistics.

Pengcheng He, Jianfeng Gao, and Weizhu Chen. 2023. DeBERTav3: Improving deBERTa using ELECTRA-style pre-training with gradient-disentangled embedding sharing. In *The Eleventh International Conference on Learning Representations*.

Pengcheng He, Xiaodong Liu, Jianfeng Gao, and Weizhu Chen. 2021. Deberta: Decoding-enhanced bert with disentangled attention. *Preprint*, arXiv:2006.03654.

Lei Huang, Weijiang Yu, Weitao Ma, Weihong Zhong, Zhangyin Feng, Haotian Wang, Qianglong Chen, Weihua Peng, Xiaocheng Feng, Bing Qin, et al. 2025. A survey on hallucination in large language models: Principles, taxonomy, challenges, and open questions. *ACM Transactions on Information Systems*, 43(2):1–55.

Ziwei Ji, Nayeon Lee, Rita Frieske, Tiezheng Yu, Dan Su, Yan Xu, Etsuko Ishii, Ye Jin Bang, Andrea Madotto, and Pascale Fung. 2023. Survey of hallucination in natural language generation. *ACM Comput. Surv.*, 55(12).

Albert Q. Jiang, Alexandre Sablayrolles, Arthur Mensch, Chris Bamford, Devendra Singh Chaplot, Diego de las Casas, Florian Bressand, Gianna Lengyel, Guillaume Lample, Lucile Saulnier, Lélio Renard Lavaud, Marie-Anne Lachaux, Pierre Stock, Teven Le Scao, Thibaut Lavril, Thomas Wang, Timothée Lacroix, and William El Sayed. 2023. Mistral 7b. *Preprint*, arXiv:2310.06825.

Albert Q. Jiang, Alexandre Sablayrolles, Antoine Roux, Arthur Mensch, Blanche Savary, Chris Bamford, Devendra Singh Chaplot, Diego de las Casas, Emma Bou Hanna, Florian Bressand, Gianna Lengyel, Guillaume Bour, Guillaume Lample, Lélio Renard Lavaud, Lucile Saulnier, Marie-Anne Lachaux, Pierre Stock, Sandeep Subramanian, Sophia Yang, Szymon Antoniak, Teven Le Scao, Théophile Gervet, Thibaut Lavril, Thomas Wang, Timothée Lacroix, and William El Sayed. 2024a. Mixtral of experts. *Preprint*, arXiv:2401.04088.

Dongfu Jiang, Yishan Li, Ge Zhang, Wenhao Huang, Bill Yuchen Lin, and Wenhu Chen. 2024b. TIGER-Score: Towards building explainable metric for all text generation tasks. *Transactions on Machine Learning Research*.

Seungone Kim, Jamin Shin, Yejin Cho, Joel Jang, Shayne Longpre, Hwaran Lee, Sangdoo Yun, Seongjin Shin, Sungdong Kim, James Thorne, and Minjoon Seo. 2024. Prometheus: Inducing fine-grained evaluation capability in language models. In *The Twelfth International Conference on Learning Representations*.

Anastasia Krithara, Anastasios Nentidis, Konstantinos Bougiatiotis, and Georgios Paliouras. 2023. Bioasq-qa: A manually curated corpus for biomedical question answering. *Scientific Data*, 10(1).

Philippe Laban, Tobias Schnabel, Paul N. Bennett, and Marti A. Hearst. 2022. SummaC: Re-visiting NLI-based models for inconsistency detection in summarization. *Transactions of the Association for Computational Linguistics*, 10:163–177.

Yang Liu, Dan Iter, Yichong Xu, Shuohang Wang, Ruochen Xu, and Chenguang Zhu. 2023a. G-eval: Nlg evaluation using gpt-4 with better human alignment. In *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing*, pages 2511–2522.

Yang Liu, Dan Iter, Yichong Xu, Shuohang Wang, Ruochen Xu, and Chenguang Zhu. 2023b. G-eval: NLG evaluation using gpt-4 with better human alignment. In *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing*, pages 2511–2522, Singapore. Association for Computational Linguistics.

Sewon Min, Kalpesh Krishna, Xinxi Lyu, Mike Lewis, Wen-tau Yih, Pang Koh, Mohit Iyyer, Luke Zettlemoyer, and Hannaneh Hajishirzi. 2023. FActScore: Fine-grained atomic evaluation of factual precision in long form text generation. In *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing*, pages 12076–12100, Singapore. Association for Computational Linguistics.

Anshuman Mishra, Dhruvesh Patel, Aparna Vijayakumar, Xiang Lorraine Li, Pavan Kapanipathi, and Kartik Talamadupula. 2021. Looking beyond sentence-level natural language inference for question answering and text summarization. In *Proceedings of the 2021 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pages 1322–1336, Online. Association for Computational Linguistics.

Guy Mor-Lan and Effi Levi. 2024. Exploring factual entailment with NLI: A news media study. In *Proceedings of the 13th Joint Conference on Lexical and Computational Semantics (*SEM 2024)*, pages 190–199, Mexico City, Mexico. Association for Computational Linguistics.

Hithesh Sankararaman, Mohammed Nasheed Yasin, Tanner Sorensen, Alessandro Di Bari, and Andreas Stolcke. 2024. Provenance: A light-weight fact-checker for retrieval augmented LLM generation output. In *Proceedings of the 2024 Conference on Empirical Methods in Natural Language Processing: Industry Track*, pages 1305–1313, Miami, Florida, US. Association for Computational Linguistics.

Thomas Scialom, Paul-Alexis Dray, Sylvain Lamprier, Benjamin Piwowarski, Jacopo Staiano, Alex Wang, Patrick Gallinari, et al. 2021. Questeval: Summarization asks for fact-based evaluation. In *Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing*, pages 6594–6604. Association for Computational Linguistics.

Damien Sileo. 2023. tasksource: Structured dataset preprocessing annotations for frictionless extreme multi-task learning and evaluation. *arXiv preprint arXiv:2301.05948*.

Julius Steen, Juri Opitz, Anette Frank, and Katja Markert. 2023. With a little push, NLI models can robustly and efficiently predict faithfulness. In *Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 2: Short Papers)*, pages 914–924, Toronto, Canada. Association for Computational Linguistics.

Gemma Team, Morgane Riviere, Shreya Pathak, Pier Giuseppe Sessa, Cassidy Hardin, Surya Bhupatiraju, Léonard Hussenot, Thomas Mesnard, Bobak Shahriari, Alexandre Ramé, Johan Ferret, Peter Liu, Pouya Tafti, Abe Friesen, Michelle Casbon, Sabela Ramos, Ravin Kumar, Charline Le Lan, Sammy Jerome, Anton Tsitsulin, Nino Vieillard, Piotr Stanczyk, Sertan Girgin, Nikola Momchev, Matt Hoffman, Shantanu Thakoor, Jean-Bastien Grill, Behnam Neyshabur, Olivier Bachem, Alanna Walton, Aliaksei Severyn, Alicia Parrish, Aliya Ahmad, Allen Hutchison, Alvin Abdagic, Amanda Carl, Amy Shen, Andy Brock, Andy Coenen, Anthony Laforge, Antonia Paterson, Ben Bastian, Bilal Piot, Bo Wu, Brandon Royal, Charlie Chen, Chintu Kumar, Chris Perry, Chris Welty, Christopher A. Choquette-Choo, Danila Sinopalnikov, David Weinberger, Dimple Vijaykumar, Dominika Rogozińska, Dustin Herbison, Elisa Bandy, Emma Wang, Eric Noland, Erica Moreira, Evan Senter, Evgenii Eltyshev, Francesco Visin, Gabriel Rasskin, Gary Wei, Glenn Cameron, Gus Martins, Hadi Hashemi, Hanna Klimczak-Plucińska, Harleen Batra, Harsh Dhand, Ivan Nardini, Jacinda Mein, Jack Zhou, James Svensson, Jeff Stanway, Jetha Chan, Jin Peng Zhou, Joana Carrasqueira, Joana Iljazi, Jocelyn Becker, Joe Fernandez, Joost van Amersfoort, Josh Gordon, Josh Lipschultz, Josh Newlan, Ju yeong Ji, Kareem Mohamed, Kartikeya Badola, Kat Black, Katie Millican, Keelin McDonell, Kelvin Nguyen, Kiranbir Sodhia, Kish Greene, Lars Lowe Sjoesund, Lauren Usui, Laurent Sifre, Lena Heuermann, Leticia Lago, Lilly McNealus, Livio Baldini Soares, Logan Kilpatrick, Lucas Dixon, Luciano Martins, Machel Reid, Manvinder Singh, Mark Iverson, Martin Görner, Mat Velloso, Mateo Wirth, Matt Davidow, Matt Miller, Matthew Rahtz, Matthew Watson, Meg Risdal, Mehran Kazemi, Michael Moynihan, Ming Zhang, Minsuk Kahng, Minwoo Park, Mofi Rahman, Mohit Khatwani, Natalie Dao, Nenshad Bardoliwalla, Nesh Devanathan, Neta Dumai, Nilay Chauhan, Oscar Wahltinez, Pankil Botarda, Parker Barnes, Paul Barham, Paul Michel, Pengchong Jin, Petko Georgiev, Phil Culliton, Pradeep Kuppala, Ramona Comanescu, Ramona Merhej, Reena Jana, Reza Ardeshir Rokni, Rishabh Agarwal, Ryan Mullins, Samaneh Saadat, Sara Mc Carthy, Sarah Perrin, Sébastien M. R. Arnold, Sebastian Krause, Shengyang Dai, Shruti Garg, Shruti Sheth, Sue Ronstrom, Susan Chan, Timothy Jordan, Ting Yu, Tom Eccles, Tom Hennigan, Tomas Kocisky, Tulsee Doshi, Vihan Jain, Vikas Yadav, Vilobh Meshram, Vishal Dharmadhikari, Warren Barkley, Wei Wei, Wenming Ye, Woohyun Han, Woosuk Kwon, Xiang Xu, Zhe Shen, Zhitao Gong, Zichuan Wei, Victor Cotruta, Phoebe Kirk, Anand Rao, Minh Giang, Ludovic Peran, Tris Warkentin, Eli Collins, Joelle Barral, Zoubin Ghahramani, Raia Hadsell, D. Sculley, Jeanine Banks, Anca Dragan, Slav Petrov, Oriol Vinyals, Jeff Dean, Demis Hassabis, Koray Kavukcuoglu, Clement Farabet, Elena Buchatskaya, Sebastian Borgeaud, Noah Fiedel, Armand Joulin, Kathleen Kenealy, Robert Dadashi, and Alek Andreev. 2024. Gemma 2: Improving open language models at a practical size. *Preprint*, arXiv:2408.00118.

Prasetya Utama, Joshua Bambrick, Nafise Sadat Moosavi, and Iryna Gurevych. 2022. Falsesum: Generating document-level nli examples for recognizing factual inconsistency in summarization. In *Proceedings of the 2022 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pages 2763–2776.

Neeraj Varshney, Wenlin Yao, Hongming Zhang, Jianshu Chen, and Dong Yu. 2023. A stitch in time saves nine: Detecting and mitigating hallucinations of llms by validating low-confidence generation. *Preprint*, arXiv:2307.03987.

Yufei Wang, Wanjun Zhong, Liangyou Li, Fei Mi, Xingshan Zeng, Wenyong Huang, Lifeng Shang, Xin Jiang, and Qun Liu. 2023. Aligning large language models with human: A survey. *arXiv preprint arXiv:2307.12966*.

Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, Ed H. Chi, Quoc Le, and Denny Zhou. 2022. Chain of thought prompting elicits reasoning in large language models. *CoRR*, abs/2201.11903.

Yue Zhang, Yafu Li, Leyang Cui, Deng Cai, Lemao Liu, Tingchen Fu, Xinting Huang, Enbo Zhao, Yu Zhang, Yulong Chen, Longyue Wang, Anh Tuan Luu, Wei Bi, Freda Shi, and Shuming Shi. 2023. Siren's song in the ai ocean: A survey on hallucination in large language models. *Preprint*, arXiv:2309.01219.

Yiran Zhao, Jinghan Zhang, I Chern, Siyang Gao, Pengfei Liu, Junxian He, et al. 2024. Felm: Benchmarking factuality evaluation of large language models. *Advances in Neural Information Processing Systems*, 36.

Lianmin Zheng, Wei-Lin Chiang, Ying Sheng, Siyuan Zhuang, Zhanghao Wu, Yonghao Zhuang, Zi Lin, Zhuohan Li, Dacheng Li, Eric Xing, et al. 2023. Judging llm-as-a-judge with mt-bench and chatbot arena. *Advances in Neural Information Processing Systems*, 36:46595–46623.

## A LLM Prompts:

The prompts used for LLM inferences on all four tasks are illustrated in Table 4.

---

**TEXT SUMMARIZATION PROMPT**

---

Summarize the given article by including the following key points:
Objective: What is the main research question or objective of the study?
Background: What is the context or rationale for the study?
Methods: What study design, population, and methodologies were used?
Key Findings: What are the most significant results or discoveries from the study?
Conclusions: What conclusions do the authors draw from their findings?
Clinical Relevance: How might the study's findings impact medical practice or patient care?
Scientific Article: article  Summary:

---

**LAY SUMMARIZATION PROMPT**

---

You will be provided a scientific article. Your task is to write a lay summary that accurately conveys the background, methods, key findings, and significance of the research in non-technical language understandable to a general audience. Guidelines for crafting a lay summary:
Craft a detailed summary that explains the research findings and their implications, providing thorough explanations where necessary.
Ensure factual accuracy and alignment with the research presented in the abstract, elaborating on key points and methodologies.
Highlight the main findings and their implications for real-world scenarios, delving into specific mechanisms or methodologies used in the study and their broader significance.
Incorporate descriptive language to explain complex concepts.
Maintain a balanced tone that is informative and engaging, avoiding technical jargon or overly formal language.
Ensure the summary provides sufficient depth and context to guide the reader through the research journey and address potential questions or areas of confusion.
Scientific Article: article
Summary:

---

**RETRIEVAL AUGMENTED GENERATION PROMPT**

---

Give a simple answer to the question based on the provided context.
QUESTION: question
CONTEXT: context

---

**OPEN-ENDED GENERATION PROMPT**

---

Give a simple answer to the question based on your best knowledge.
QUESTION: question

---

Table 4: The prompt in the Benchmark for LLM generation output for all tasks.

## B Technical Details

### B.1 LLM Generations

The inference procedure was done Together AI Inference[4]. We used the instruct-tuned or chat versions of the models. As for GPT-4o mini, we used the OpenAI API and the latest snapshot available, `gpt-4o-mini` from Sep 13th, 2024. The checkpoints used for LLM inferences of the open-source models using Together AI are summarized in Table 5.

| Model | checkpoint |
|---|---|
| Llama 3.1 8b | meta-llama/Meta-Llama-3.1-8B-Instruct-Turbo |
| Llama 3.1 70b | meta-llama/Meta-Llama-3.1-70B-Instruct-Turbo |
| Mistral 7b | mistralai/Mistral-7B-Instruct-v0.3 |
| Mixtral 8x7b | mistralai/Mixtral-8x7B-Instruct-v0.1 |
| Gemma 2 9b | google/gemma-2-9b-it |

Table 5: Together AI checkpoints of all LLMs that were used during Inferences.

### B.2 Benchmark Computations

We used the OpenAI API[5] and the latest snapshot available, `GPT-4o-mini` from Sep 13th, 2024 for Fact-checking using Chain-of-Thought prompting. We leveraged Nvidia V100-16GB and Nvidia A100-80GB GPUs for performing fact-checking.

---

[4] https://www.together.ai/

[5] https://platform.openai.com/

# LLMs Information Flow Diagnostic:
# Memory-based Evidence from Random Matrix Theory

**Sami Diaf**

Department of Socioeconomics
Universität Hamburg
`sami.diaf@uni-hamburg.de`

## Abstract

The diagnostic of neural networks, particularly Large Language Models (LLMs), remains a critical aspect of today's AI-powered solutions, whose training data are not available to users for testing purposes. Practitioners usually aim to fine-tune their models to maximize the accuracy, by leveraging the traditional test metrics, whose application on large models remains expensive. Recent advances considered layer-based norms and power-law metrics for a robust meta-analysis, without the need to access training and test data. Inherently, elements from Random Matrix Theory were used to reveal inner correlation patterns and size scales within each layer, so to detect bottlenecks in pre-trained models. This article extends the use of such schemes by analyzing memory dynamics and the probabilistic properties of power-law metrics to study the information flow within specific LLMs. Taken on a pretained German LLM (*LLaMmlein*) and its original English model (*TinyLlama*), this approach confirmed embedded self-similar, fractal properties of power-law metrics, hinting heavy tails and long-range correlations in the training process with a substantial amount of undertrained layers. This variability was found to be slightly persistent in the original English TinyLlama model and its German version, however the latter's chat version exhibits a pure randomness in its metrics. Findings stress out the role of attention mechanism as the main driver of LLMs training issues, while language-specific structures may cause metrics' distortions, hence altering the inter-layer information transmission as a component of the training process.

## 1 Introduction

The advent of neural networks, coupled with intensive computational innovations, popularized the use of deep learning as a modeling standard, outperforming other existing machine learning algorithms. Although the widespread use of such capabilities opened new research areas, deep neural networks (DNNs) remain black box models, whose effectiveness depends on complex hyperparameter optimization (Wu et al., 2019) to achieve a robust training. This forced practitioners to adopt expensive feature engineering schemes, without clearly setting up a strong theoretical background for users (Martin et al., 2021).

Large Language Models (LLMs) have been extensively designed, as large scale models, to accomplish several complex tasks in Natural Language Processing (NLP). Tuning and testing such models require extensive learning time (Burns et al., 2025), while training and test data are not always publicly available. Moreover, such DNNs are based on transformers (Vaswani et al., 2017) and require a special attention because they feature memory mechanisms, as for multihead attention and BiLSTM (Graves and Schmidhuber, 2005). Although these memory-based architectures are complex to handle, they became the default choice for many NLP architectures, as for the popular BERT model (Devlin et al., 2019).

The term *memory* refers, for the particular case of DNNs, to any mechanism by which a model or agent stores, retrieves and uses historical information (Zhang et al., 2024b), whether internally or externally. This paper considers the memory stemming from the information exchanged between layers, that is the output flow of each layer in the architecture, given by its weight matrix.

Random Matrix Theory (RMT) (Tulino and Verdú, 2004) is considered as the central limit theorem for matrix analysis and was used to study the overall performance of DNNs (Martin and Mahoney, 2021), on the basis of extracted eigenvalues of each weight matrix in the architecture. While earlier approaches considered mapping neural networks to a Gaussian process (Jacot et al., 2018), Martin et al. (2021) set up a practical background to identify similarities in the learning process of multiple DNNs, particularly fitting issues and the

*bona fide* of different regularization schemes to reduce correlations inside each layer. This extended the concept of Self-Regularization theory (Malevergne and Sornette, 2004), which assumes the generic existence of a self-organized macroscopic state in any large multivariate system. Martin et al. (2021) came to the conclusion that an implicit self-regularization at DNNs was prevailing, at the contrast of explicit regularization (L1 and L2) constraining the norm of weight matrices.

This new field of research set up effective generalization metrics detailing the inner functioning of DNNs, especially the learning process, the inter-layer information flow and the intra-layer asymptotic convergence (Martin et al., 2021). It borrows elements from statistical mechanics and was used for many applications as for cyber threat detection (Ferrag et al., 2024) and the description of feature learning applications (Seroussi et al., 2023).

In parallel to the use of power laws (PL) in various scientific fields, pattern similarities were studied under the name of *fractal analysis*, defining the behavior of self-similar patterns whose occurrence is not purely random, but follows a power-law behavior (Mandelbrot, 1982). The *fractality* is an essential feature in language theory, denoting the complexity stemming from word usage (Hiver et al., 2022), and was recently used in information processing (Wang et al., 2024). It fits the study of the information correlation proposed by Martin et al. (2021) which relies on a power-law fit over heavy-tailed distributions.

While the training quality of popular NLP and Computer Vision models came to scrutiny via norms and PL-based metrics (Yang et al., 2023), it ignored their inter-layer information exchange as a component of the training process. This concern is particularly determinant for LLMs, whose complex architecture features two distinct types of attention mechanisms (Vaswani et al., 2017; Martin et al., 2021), as a key component a transformer.

Thus, this paper enriches the existing DNNs empirical methodology by investigating the existence of pattern similarity in the information transmission on selected LLMs trained over English and German corpora. It extends the layer-based meta-analysis on such big architectures and details inter-layer persistence behavior. The latter reveals short/long term variations in the training process, whose non-linearity is linked to underfitted layers.

For this aim, two German LLMs, namely *LLaMmlein_1B* model[1] and a lightweight, small-scale version *LLaMmlein_120M* model[2], were used in this paper to conduct a transfer learning experiment, along the English *TinyLlama*, who served in training the *LLaMmlein*.

Aside from a meta-analysis on each selected LLM following Martin et al. (2021), an additional memory check was conducted to dissect hidden trends in the PL-based metrics. It revealed mild persistency and underfitting of metrics featuring information correlation and the size scale. Metrics based solely on information correlation were found to indicate heavy-tailed distribution of the eigenvalues and a high persistence, denoting the importance of the size scale in the information flow analysis.

Findings indicate layers exhibit substantial underfitting properties in both languages, mainly due to attention mechanisms. Original TinyLlama (Zhang et al., 2024a), both the full and the chat versions, have a mild persistent flow of information, compared to the German LLaMmlein whose lightweight version is though slightly anti-persistent. The size scale, measured by the maximum eigenvalue, proved to be important in harmonizing the per-layer metrics. Differences in results obtained from English and German LLMs could be explained by the morphologically-rich characteristic of the German language, known to be a SOV (Subject-Object-Verb), while English language exhibits a less complex SVO structure (Vikner, 2019).

The paper outlines the use of Random Matrix Theory in DNNs analysis (Section 2), then details the Rescaled Range Analysis (Hurst, 1951), as a method to study fractal properties and persistency measurement (Section 3). Section 4 features two language-based applications on English and German LLMs and compares their metrics and persistency measurements.

## 2 Random Matrix Theory

Train and test data have been the de facto tools to assess machine learning models in general, and neural networks in particular. In the absence of such data, elements from Random Matrix Theory were applied on final weight matrices of neural networks (Martin and Mahoney, 2021) to check their asymptotic convergence. It resulted several norms and metrics, whose statistical properties were found to

---

[1] https://huggingface.co/LSX-UniWue/LLaMmlein_1B

[2] https://huggingface.co/LSX-UniWue/LLaMmlein_120M

match DNNs accuracy, without accessing data used to train the models (Martin et al., 2021). In other terms, this strategy permits to discover whether a layer learned too much from the noise (overfitting) or alternatively has not learned enough from the signal (underfitting), assuming data stem from two components: signal and noise.

The *WeightWatcher* open source tool (Martin et al., 2021) investigates the weight matrix $W$ of a given DNN layer, by analyzing its spectral properties. While every element of the weight matrix[3] $W_{ij}$ is assumed to follow a normal distribution $\mathcal{N}(0, \sigma^2)$, the empirical correlation (Wishart) matrix $\boldsymbol{X} = \frac{1}{N} W^\intercal W$ is taken as the basis for quality assessment, by extracting its eigenvalues spectrum.

The Marchenko-Pastur (MP) distribution (Marchenko and Pastur, 1967) considers the spectrum of eigenvalues bounded between $\lambda_-$ and $\lambda_+$ as relevant to the noise randomness. Its probability density $f(\lambda)$ is given for a $T \times N$ matrix and a noise level $\sigma^2$ as:

$$f(\lambda) = \begin{cases} \frac{N}{T} \frac{\sqrt{(\lambda_+ - \lambda)(\lambda - \lambda_-)}}{2\pi\sigma^2} & \text{if } \lambda \in [\lambda_-, \lambda_+], \\ 0 & \text{if } \lambda \notin [\lambda_-, \lambda_+]. \end{cases}$$

where $\lambda_- = \sigma^2(1 - \sqrt{\frac{T}{N}})^2$ and $\lambda_+ = \sigma^2(1 + \sqrt{\frac{T}{N}})^2$

The eigenvalues distribution, plotted as a histogram using the Empirical Spectral Density (ESD), is an informative feature of the randomness prevailing in every layer constituting the DNN, in addition to reveal inter-layer differences.

Because many matrices hold strongly correlated elements, the MP distribution is used to empirically evaluate a noisy spectrum of eigenvalues, that could be separated from other eigenvalues representing the signal.

Martin and Mahoney (2021) found most weight matrices in DNNs exhibit heavy-tailed distributions of eigenvalues as they become increasingly correlated, suggesting rather drawing elements from power-law generated data, as for Pareto distribution. This concept, known as Heavy-Tailed Self-Regularization (HT-SR) theory, is linked to situations where separating the noise from the signal becomes difficult to achieve, as eigenvalues are in this case better modeled via heavy-tailed distribu-

---

[3]A layer with multiple weight matrices will have a single concatenated weight matrix (Martin et al., 2021).

tions (Malevergne and Sornette, 2004), rather than a simple MP distribution.

For this aim, Martin and Mahoney (2021) estimated a truncated power-law fit (Clauset et al., 2009) over the MP curve, yielding the exponent $\alpha$ from the equation ESD–eigenvalues: $\rho(\lambda) \sim \lambda^{-\alpha}$ for $\lambda \in [\lambda_-, \lambda_+]$. The amplitude of the PL-exponent $\alpha$ is considered as the *information correlation* index within each weight matrix, denoting the strength of the existing element-wise correlations. Moreover, the $\alpha$ exponent is indeed a power-law fit that can be considered as a complexity index or a *fractal dimension* (Mandelbrot, 1982).

Based on the eigenvalues spectrum $\lambda_i$ of each correlation matrix $\boldsymbol{X}$, several metrics were used as for:

- Frobenius norm : $\|W\|_F^2 = \|X\|_F = \sum_{i=1}^{M} \lambda_i^2$

- Spectral norm : $\|W\|_\infty = \|X\|_\infty = \lambda_{max}$

- Weighted $\alpha$ : $\hat{\alpha} = \alpha Log \lambda_{max}$

- $\alpha$ norm (Shatten-norm) : $\|W\|_{2\alpha}^{2\alpha} = \|X\|_\alpha^\alpha = \sum_{i=1}^{M} \lambda_i^\alpha$

where $\lambda_i$ is the $i^{th}$ eigenvalue of $\boldsymbol{X}$, $\lambda_{max}$ is the maximum eigenvalue and $\alpha$ is the fitted power-law exponent, usually truncated because it needs defining specific lower and upper bounds, respectively $\lambda_-$ and $\lambda_+$. For instance, Figure 1 reports simulations yielding random-like eigenvalues fitted with a scale-invariant Marchenko-Pastur curve between $\lambda_- \simeq 0.31$ and $\lambda_+ \simeq 1.17$ and spikes (signal) associated with $\lambda_i > \lambda_+$. The PL-fit yields a value of 0.571 for $\alpha$.

The plain $\alpha$ metric is a scale-invariant, weak estimation of the information correlation, as it ignores the size scale ($\lambda_{max}$) within each layer. The latter remains an important determinant of HT-SR because DNNs are known to be non-linear, while LLMs particularly feature attention layers with large matrices. For small values of $\alpha$, the size scale $\lambda_{max}$ was found to be a good proxy for estimating the difference between the noise and the signal, however, for higher values of $\alpha$ (HT-SR), the signal gets mixed with the noise and $\lambda_{max}$ is non-informative.

A clear distinction between norm-metrics and PL-based metric was given when studying the performance of several DNNs models (Martin et al., 2021; Yang et al., 2023). They concluded that

Figure 1: Marchenko-Pastur distribution simulated 1,000 times on the correlation matrix of an initial random matrix with $\frac{T}{N}$=10 and $\sigma^2 = \frac{2}{3}$ . $\alpha$ is the PL-exponent of the Marchenko-Pastur fit over the interval $[\lambda_-,\lambda_+]$.

PL-based metrics, aside from being good proxies for overall accuracy measurements, remain robust in detecting potential bottlenecks and training issues than norm-based metrics. Hence, PL exponent remains a robust empirical metric to asses well-trained DNNs and quantify the layer-wise correlation flow (Martin et al., 2021).

In practice, $\alpha$ was found to match an ideal DNN fit when approaching 2. This means the DNN model performs well as it facilitates the propagation of information/features across layers, because it learns from both data signal and noise. Values in the interval [4,6] are proxies of underfitting situations (not learning enough from the signal), while lower values equaling 1.5 are synonyms of overfitting (learning too much from the noise) (Martin et al., 2021). Large values of $\alpha > 6$ are associated with a pure randomness, which requires the aspect ratio $\frac{T}{N}$ to differentiate layers.

Because the size of DNNs layers changes according to adopted architectures, Martin et al. (2021) proposed to weight the $\alpha$ with the size scale to produce the weighted $\alpha$ metric. It was found that for small values, the weighted $\alpha$ approximates well the $\alpha$ Shatten-norm; the latter weighs the $\alpha$ exponent for all eigenvalues within the layer.

Martin et al. (2021) reported that weighted $\alpha$ and log $\alpha$ norm correlate at a higher level for well trained models. The size scale, given by $\lambda_{max}$, could be informally linked to situations where input clusters are at a greater distance. This means the size scale is related, in the case of LLMs, to the language morphologic aspects (sentence structures).

Particularly in LLMs, distortions in the series of PL exponents is called *scale collapse*, mostly linked to transformers (Vaswani et al., 2017; Lefaudeux et al., 2022). As memory-based blocks of layers, transformers feature a complex inner structure usually yielding larger weight matrices.

The study of such variations and the training process requires detailing the information flow throughout the whole network. The adoption of advanced tool for self-similar patterns, known as fractals (Mandelbrot, 1982) is clearly indicated to test the persistency hypothesis on trained DNNs. Persistent behavior of the aforementioned metrics reinforces the hypothesis of a strong, correlated inter-layer linkage propping up the information flow. One can assert that anti-persistency of PL-metrics may indicate colliding trends that alter the training process and the inter-layer dynamics, while persistency may reinforce the hypothesis of a harmonized network design that better captures long-range dependencies via attention layers.

## 3 Fractal Analysis

Mandelbrot tried first to uncover repeated patterns able to explain the randomness of irregular shapes (Mandelbrot, 1982), as exemplified by Koch's snowflake. This led to the concept of self-similar patterns, which stands for scale-dependent shapes with a known geometry. Hence, the *fractal* analysis was first established as a research field in geometry having a wide range of applications, from physics to hydrology. The fractal theory relies on the definition of a fractal dimension, a hidden variable that quantifies the irregularity of shapes found in many objects.

In time series analysis, the fractal approach was first featured when studying the Nile river flooding history. Hurst (1951) designed the *Rescaled Range (R/S) Analysis* and reckoned the Hurst exponent as a measure of a time series memory, later corrected by Mandelbrot and extended to the fractional Brownian motion (Mandelbrot and van Ness, 1968) when studying cotton prices in the United States.

The R/S algorithm takes the variations of a given time series of length $T$ and divides them into $N$ adjacent intervals of length $\tau$, where $T = N\tau$. For each interval, the average value is computed and a new time series is created as accumulated deviations from the arithmetic mean values (hereafter named profile). The difference (range) between the

maximum and the minimum value of the profile, and the standard deviation of the original time series for each interval, are calculated. Each range is standardized by the corresponding standard deviation and forms a rescaled range so that the average *rescaled range* for a given interval of length $(R/S)_\tau$ is calculated.

The rescaled range scales are given by $(R/S)_\tau \approx c\tau^H$, where $c$ is a finite constant independent of $\tau$ (Taqqu et al., 1995). To estimate the power law relationship, a simple log-log ordinary least squares regression is used for: $\log (R/S)_\tau \approx \log c + H \times \log \tau$, where $H$ is the estimated Hurst exponent (Barunik and Kristoufek, 2010). $R/S$ analysis was shown to be biased for small $\tau$ (Couillard and Davison, 2005), and empirical application considered rather the expected Hurst exponent (Weron, 2011). Values of $H$ exceeding 0.5 are proxies of a persistent behavior resulting from long-range correlations, while values less than 0.5 are anti-persistent. A Hurst exponent not significantly different from 0.5 is associated to the standard Brownian motion. The Hurst exponent $H$ is also a proxy of the fractal dimension $D$ in time series, linked by the relationship: $D = 2 - H$.

Given the relatively reduced number of layers in most DNNs, this article considers the existence of a single fractal dimension, approached by the Hurst exponent. For each layer in an LLM, PL-metrics are computed on the related weight matrix, yielding three different series across the whole LLM to run the R/S Analysis on each one of them.

# 4 Application

The study of memory properties of specific LLMs is conducted on the weight matrices, stored after achieving the LLMs training. PL-based metrics adopted by Martin et al. (2021) were previously found to be robust when assessing hundreds of LLMs, outperforming simple algebraic norms (Frobenius and spectral norms).

The weighted $\alpha$ and log $\alpha$ norm are compound metrics computed from a truncated PL-fit of the eigenvalues and the size scale. These two metrics will have a particular attention in this section, as they go in-line with the PL-exponent yielded by the R/S Analysis, known as the Hurst exponent. The purpose lies on investigating the inter-layer dynamic flow using above two metrics and uncover potential variability known as *scale collapse* (Martin et al., 2021), which is assumed to reveal dys-

functions in the learning process. The $\alpha$ series will not be considered for the R/S analysis, as it ignores the size scale.

The selected LLMs are publicly available and their PyTorch versions (Paszke et al., 2019) were used to run the *WeightWatcher* diagnostic tool. The R/S analysis was performed on the basis of estimated PL-metrics, whose relatively reduced size requires a corrected version of the Hurst exponent (Weron, 2011) reported in Table 2.

## 4.1 English *TinyLlama*

TinyLlama model (Zhang et al., 2024a) was trained on a complex architecture featuring flash attention 2 and various fused schemes, comprising xFormers (Lefaudeux et al., 2022) as a research tool for accelerated transformers.

Figure 2 displays the per-layer metrics for the TinyLlama 1.1B model trained over 155 layers. The weighted $\alpha$ and the log $\alpha$ norm are highly correlated and clearly separable from the simple $\alpha$ metric, which exhibits a pronounced variability. This denotes the importance of the size scale, absent from the $\alpha$ metric, but present in the two others. Similar patterns were found in the TinyLlama 1.1B chat model (Figure 3), although its first layers are less pronounced then the original model.

The variability of the above metrics is a result of heavy-tailed eigenvalues distributions associated to a *scale collapse*. This denotes implicit changes or perturbations that occurred when training the model, likely due to distillation, data augmentation or fine-tuning.

Both LLMs feature a relatively high number of layers found to be under-trained, as reported in Table 1. These demonstrate high $\alpha$ values and are linked to value-type (V) self attention layers (having a rank of 256). They are particularly aggregated representations of the words in context (Vaswani et al., 2017), compared to query (Q) and key (K) matrices. The relative low number of over-trained layers confirms difficulties of fine tuning LLMs who are over-trained (Springer et al., 2025).

First layers, usually associated with higher metrics due to their effective normalization (Martin et al., 2021), do not exhibit here higher values of weighted $\alpha$ and the log $\alpha$ norm, compared to what was reported in Martin et al. (2021).

Table 2 reveals a slight persistency of the weighted $\alpha$ and log $\alpha$ norm metrics for the LLM chat version (Hurst exponent respectively 0.60

Figure 2: PL metrics estimated from TinyLlama 1.1B model



Figure 3: PL metrics estimated from TinyLlama 1.1B Chat

and 0.61), while the full model exhibits a non-persistent, Brownian-like behavior (Hurst exponent 0.51 each). The buildup of the chat version proved to have more inter-layer information than the original model, as a result of intensive fine-tuning on synthetic dialogues provided by Zephyr (Tunstall et al., 2023).

Both LLMs show similar PL-metric patterns and persistence, reinforcing the hypothesis of a strong transfer learning between the original model TinyLlama 1.1B and its chat version. The metric correlations of weighted $\alpha$ and log $\alpha$ norm are almost identical, respectively 0.879 and 0.887.

## 4.2 German *LLaMmlein*

The layer-to-layer information flow, as given by three metrics in Figure 4 and Figure 5, demonstrates key differences between the German LLaMmlein and its lightweight version (LLaMmlein 120M chat). The latter features 85 layers, compared to the 155 comprised in the former. Weighted $\alpha$ and log $\alpha$ norm are highly correlated in both models, however, the lightweight version displays a relatively stable $\alpha$ metric, not as variable as in the LLaMmlein 1B model, whose metrics have long-range correlations (Hurst exponent 0.61 in Table 2.

Higher values of $\alpha$ for LLaMmlein 1B are associated with V self attention layers of rank 256 (Figure 4), that carry context-based information of each sentence/word fed to the LLM. The lightweight version (LLaMmlein 120M) presents the lowest rate of under-trained layers, despite its reduced depth. This means this abridged version does not suffer from over-parametrization, relative to the amount of data. However, slight differences in the Hurst exponent values indicate a weak anti-persistency of the weighted $\alpha$ (Hurst exponent 0.46) compared to Brownian-like log $\alpha$ norm (Hurst exponent 0.52).

The impact of the size scale ($\lambda_{max}$) seems to be mild in the lightweight version, in comparison with the full model. This explains why the information correlation series $\alpha$ does not feature very high values in the lightweight model and exhibit a relative stability compared to the full model. The size scale has, particularly for the lightweight version, a linguistic feature embedded in the dataset[4].

The German language features a SOV structure (Vikner, 2019), at the contrary of the common SVO structures found in English and French. This considers German as a morphologically-rich language (Günther et al., 2019) whose structure is complex but rich, compared to English. Moreover, German LLMs are mostly trained on the basis of existing English and/or Multilingual LLMs, while recent attempts proposed a data curation methodology to improve LLMs training (Burns et al., 2025).



Figure 4: PL metrics estimated from LLaMmlein 1B model.

## 5 Conclusion

Machine learning models have long been associated with the train/test paradigm and the related metrics to perform quality control checks. For DNNs, practitioners use models without access

---

[4]Training data were de-duplicated on the paragraph level and filtered using a token-to-word ratio.

Figure 5: PL metrics estimated from LLaMmlein 120M model.

| Model | Overtrained | Undertrained |
|---|---|---|
| TinyLlama 1.1 B | 1.3% | 26.3% |
| TinyLlama 1.1 B Chat | 1.3% | 29.5% |
| LLaMmlein 1B | 2.9% | 28.8% |
| LLaMmlein 120M | 2.3% | 13.9% |

Table 1: Percentages of over-/under-trained layers, based on estimated $\alpha$ values, obtained from *WeightWatcher* tool (Martin and Mahoney, 2021)

to training data and are not able to perform independent accuracy tests. Elements from statistical mechanics were used to check the robustness of DNNs on the basis of their weight matrices, as information-carriers of the learning process. The use of Random Matrix Theory helped revealing embedded, heavy-tailed properties of eigenvalues via a truncated power-law fit, whose exponent is taken as a proxy of underfitting or overfitting presence in the related layer. Hybrid metrics combining power-law exponents and size scale proved to be accurate in estimating the between/within layer information flow, particularly in the case of LLMs who feature attention layers as memory-driver mechanisms. The inter-layer information flow, as an element of the training process, was found to exhibit a noticeable persistence in terms of long-range correlations. Such findings confirm the fractality of LLMs learning process and the importance of language-properties carried by data, whose complexity flags substantial underfitting issues affecting attention layers. The self-similarity analysis provides tools to detect potential training bottlenecks, but also a powerful way to assess transfer learning strategies when designing lightweight and task- and language-specific models. This proved particularly effective for the German language, whose morphologically-rich properties make the training difficult and require a special hyperparameter tuning and data processing.

| Model | $\alpha$ | Weighted $\alpha$ | Log $\alpha$ norm |
|---|---|---|---|
| TinyLlama 1.1 B | 0.63 | 0.51 | 0.51 |
| TinyLlama 1.1 B Chat | 0.49 | 0.60 | 0.61 |
| LLaMmlein 1B | 0.79 | 0.61 | 0.61 |
| LLaMmlein 120M | 0.74 | 0.46 | 0.52 |

Table 2: Estimates of Hurst exponents for each model, based on estimated $\alpha$, weighted $\alpha$ and log $\alpha$ norm, obtained from *WeightWatcher* tool (Martin and Mahoney, 2021)

# References

Jozef Barunik and Ladislav Kristoufek. 2010. On hurst exponent estimation under heavy-tailed distributions. *Physica A: Statistical Mechanics and its Applications*, 389(18):3844–3855.

Thomas F Burns, Letitia Parcalabescu, Stephan Wäldchen, Michael Barlow, Gregor Ziegltrum, Volker Stampa, Bastian Harren, and Björn Deiseroth. 2025. Aleph-alpha-germanweb: Improving german-language llm pre-training with model-based data curation and synthetic data generation.

Aaron Clauset, Cosma Rohilla Shalizi, and M. E. J. Newman. 2009. Power-law distributions in empirical data. *SIAM Review*, 51(4):661–703.

Michel Couillard and Matt Davison. 2005. A comment on measuring the Hurst exponent of financial time series. *Physica A: Statistical Mechanics and its Applications*, 348(C):404–418.

Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. 2019. Bert: Pre-training of deep bidirectional transformers for language understanding. In *Proceedings of the 2019 conference of the North American chapter of the association for computational linguistics: human language technologies, volume 1 (long and short papers)*, pages 4171–4186.

Mohamed Amine Ferrag, Mthandazo Ndhlovu, Norbert Tihanyi, Lucas C. Cordeiro, Merouane Debbah, Thierry Lestable, and Narinderjit Singh Thandi. 2024. Revolutionizing cyber threat detection with large language models: A privacy-preserving bert-based lightweight model for iot/iiot devices. *IEEE Access*, 12:23733–23750.

Alex Graves and Jürgen Schmidhuber. 2005. Framewise phoneme classification with bidirectional lstm and other neural network architectures. *Neural Networks*, 18(5):602–610. IJCNN 2005.

Fritz Günther, Eva Smolka, and Marco Marelli. 2019. 'understanding' differs between english and german: Capturing systematic language differences of complex words. *Cortex*, 116:168–175. Structure in words: the present and future of morphological processing in a multidisciplinary perspective.

Phil Hiver, Ali H. Al-Hoorie, and Reid Evans. 2022. Complex dynamic systems theory in language learning: A scoping review of 25 years of research. *Studies in Second Language Acquisition*, 44(4):913–941.

Harold Edwin Hurst. 1951. Long-term storage capacity of reservoirs. *Transactions of the American society of civil engineers*, 116(1):770–799.

Arthur Jacot, Franck Gabriel, and Clément Hongler. 2018. Neural tangent kernel: convergence and generalization in neural networks. In *Proceedings of the 32nd International Conference on Neural Information Processing Systems*, NIPS'18, page 8580–8589, Red Hook, NY, USA. Curran Associates Inc.

Benjamin Lefaudeux, Francisco Massa, Diana Liskovich, Wenhan Xiong, Vittorio Caggiano, Sean Naren, Min Xu, Jieru Hu, Marta Tintore, Susan Zhang, Patrick Labatut, Daniel Haziza, Luca Wehrstedt, Jeremy Reizenstein, and Grigory Sizov. 2022. xformers: A modular and hackable transformer modelling library. https://github.com/facebookresearch/xformers.

Y. Malevergne and D. Sornette. 2004. Collective origin of the coexistence of apparent random matrix theory noise and of factors in large sample correlation matrices. *Physica A: Statistical Mechanics and its Applications*, 331(3):660–668.

Benoit B. Mandelbrot and John W. van Ness. 1968. Fractional Brownian motions, fractional noises and applications. *SIAM Review*, 10:422–437.

Benoît B. Mandelbrot. 1982. *The fractal geometry of nature*. W. H. Freeman and Comp., New York.

V. A. Marchenko and L. A. Pastur. 1967. Distribution of eigenvalues for some sets of random matrices. *Math. USSR-Sb*, 1:422–437.

Charles H. Martin and Michael W. Mahoney. 2021. Implicit self-regularization in deep neural networks: evidence from random matrix theory and implications for learning. *Journal of Machine Learning Research*, 22(1).

Charles H. Martin, Tongsu Peng, and Michael W. Mahoney. 2021. Predicting trends in the quality of state-of-the-art neural networks without access to training or testing data. *Nature Communications*, 12(1).

Adam Paszke, Sam Gross, Francisco Massa, Adam Lerer, James Bradbury, Gregory Chanan, Trevor Killeen, Zeming Lin, Natalia Gimelshein, Luca Antiga, Alban Desmaison, Andreas Köpf, Edward Yang, Zach DeVito, Martin Raison, Alykhan Tejani, Sasank Chilamkurthy, Benoit Steiner, Lu Fang, Junjie Bai, and Soumith Chintala. 2019. Pytorch: An imperative style, high-performance deep learning library.

Inbar Seroussi, Gadi Naveh, and Zohar Ringel. 2023. Separation of scales and a thermodynamic description of feature learning in some cnns. *Nature Communications*, 14(1):908.

Jacob Mitchell Springer, Sachin Goyal, Kaiyue Wen, Tanishq Kumar, Xiang Yue, Sadhika Malladi, Graham Neubig, and Aditi Raghunathan. 2025. Overtrained language models are harder to fine-tune.

Murad S. Taqqu, Vadim Teverovsky, and Walter Willinger. 1995. Estimators for long-range dependence: An empirical study. *Fractals*, 03(04):785–798.

Antonia Tulino and Sergio Verdú. 2004. *Random Matrix Theory and Wireless Communications*. Now Foundations and Trends.

Lewis Tunstall, Edward Beeching, Nathan Lambert, Nazneen Rajani, Kashif Rasul, Younes Belkada, Shengyi Huang, Leandro von Werra, Clémentine Fourrier, Nathan Habib, Nathan Sarrazin, Omar Sanseviero, Alexander M. Rush, and Thomas Wolf. 2023. Zephyr: Direct distillation of lm alignment.

Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N. Gomez, Łukasz Kaiser, and Illia Polosukhin. 2017. Attention is all you need. In *Proceedings of the 31st International Conference on Neural Information Processing Systems*, NIPS'17, page 6000–6010, Red Hook, NY, USA. Curran Associates Inc.

Sten Vikner. 2019. *Why German is not an SVO-language but an SOV-language with V2*. AU Library Scholarly Publishing Services.

Zhenhua Wang, Fuqian Zhang, Ming Ren, and Dong Gao. 2024. A new multifractal-based deep learning model for text mining. *Information Processing and Management*, 61(1):103561.

Rafal Weron. 2011. HURST: MATLAB function to compute the Hurst exponent using R/S Analysis. HSC Software, Hugo Steinhaus Center, Wroclaw University of Science and Technology.

Jia Wu, Xiu-Yun Chen, Hao Zhang, Li-Dong Xiong, Hang Lei, and Si-Hao Deng. 2019. Hyperparameter optimization for machine learning models based on bayesian optimization. *Journal of Electronic Science and Technology*, 17(1):26–40.

Yaoqing Yang, Ryan Theisen, Liam Hodgkinson, Joseph E. Gonzalez, Kannan Ramchandran, Charles H. Martin, and Michael W. Mahoney. 2023. Test accuracy vs. generalization gap: Model selection in nlp without accessing training or testing data. In *Proceedings of the 29th ACM SIGKDD Conference on Knowledge Discovery and Data Mining*, KDD '23, page 3011–3021, New York, NY, USA. Association for Computing Machinery.

Peiyuan Zhang, Guangtao Zeng, Tianduo Wang, and Wei Lu. 2024a. Tinyllama: An open-source small language model.

Zeyu Zhang, Xiaohe Bo, Chen Ma, Rui Li, Xu Chen, Quanyu Dai, Jieming Zhu, Zhenhua Dong, and Ji-Rong Wen. 2024b. A survey on the memory mechanism of large language model based agents.

# CUPE: Contextless Universal Phoneme Encoder for Language-Agnostic Speech Processing

**Abdul Rehman**  **Jian-Jun Zhang**  **Xiaosong Yang**

Bournemouth University

Bournemouth, United Kingdom

`arehman, jjunzhang, xyang@bournemouth.ac.uk`

## Abstract

Universal phoneme recognition typically requires analyzing long speech segments and language-specific patterns. Many speech processing tasks require pure phoneme representations free from contextual influence, which motivated our development of CUPE - a lightweight model that captures key phoneme features in just 120 milliseconds, about one phoneme's length. CUPE processes short, fixed-width windows independently and, despite fewer parameters than current approaches, achieves competitive cross-lingual performance by learning fundamental acoustic patterns common to all languages. Our extensive evaluation through supervised and self-supervised training on diverse languages, including zero-shot tests on the UCLA Phonetic Corpus, demonstrates strong cross-lingual generalization and reveals that effective universal speech processing is possible through modeling basic acoustic patterns within phoneme-length windows.

## 1 Introduction

Current speech processing systems depend heavily on contextual information, creating a double-edged sword for certain tasks. While extensive context provides crucial bias toward appropriate attention mechanisms, it simultaneously makes it nearly impossible to isolate individual speech units—particularly allophones—from their contextual embeddings. Modern systems such as derivatives of wav2vec 2.0 (Baevski et al., 2020) typically analyze 300-2500ms of speech, incorporating extensive language-specific patterns and contextual dependencies. While effective for automatic speech recognition, this approach entangles phonetic content with contextual information, making it extremely difficult to disentangle the acoustic properties that define individual speech sounds.

The necessity for contextless processing emerges from two critical considerations: alignment precision and representation purity. Extended temporal windows (e.g., 500ms) reduce inter-frame discriminability as individual frame representations become increasingly influenced by surrounding context. Optimal alignment performance requires maximally discriminative frame-level representations, where each frame maintains distinct characteristics. As context window length increases, the transformer's attention mechanism progressively attenuates frame-specific features through contextual averaging, resulting in diminished temporal resolution.

For paralinguistic tasks, contextless models function as quantization preprocessing stages. When frame-level embeddings encode predominantly contextual rather than local information, this homogenization undermines the model's capacity to capture subtle local acoustic variations essential for allophone analysis and speaker-specific phonetic characterization.

Our empirical results directly challenge the assumption that more context is always better—models using 120ms of speech windows actually perform on-par if not better than those using full word context across multiple evaluation scenarios, while simultaneously providing access to pure phonemic representations less contaminated by contextual dependencies.

Our work makes three key contributions. First, we demonstrate that universal phoneme recognition can be achieved effectively with just 120ms of context, a fraction of the 300-2500ms typically used in current approaches. Second, we introduce CUPE, a lightweight architecture (30M parameters) that achieves competitive performance through focused local feature extraction. Third, we provide a feature extraction method that captures pure phonemic representations by eliminating contextual dependencies, leading to cleaner and more interpretable phoneme embeddings across languages. By operating on brief windows—approximately the duration

of a typical phoneme (Crystal and House, 1988), CUPE learns language-agnostic acoustic features that characterize phonemes universally. This focus on fundamental acoustic patterns, independent of language-specific context, enables robust cross-lingual generalization and, crucially, provides access to clean allophonic representations that are essential for understanding speaker-specific phonetic variations.

The contextless nature of our approach enables several practical applications:

- **Timestamps Alignment**: Generating time-aligned transcripts from raw text and audio. This task is critical for training downstream text-to-speech models. Since this is the main application for phoneme recognition, it helps to have as little context information in each frame so that there is a sharper contrast between frames for precise boundary detection.

- **Speech style learning**: It serves as a foundational allophone encoder. Embeddings of each frame can be used to generate acoustically pure allophone variants of base phonemes. This is also useful for training downstream text-to-speech tasks which currently rely on IPA dictionaries or sub-word tokens.

- **Robust phoneme verification**: Complementing traditional ASR systems by detecting and correcting errors that arise from over-reliance on language context.

- **Cross-linguistic research**: Generating language-agnostic phoneme representations that facilitate multilingual studies and enable more accurate speech disorder diagnostics.

Through extensive evaluation, we validate CUPE (Contextless Universal Phoneme Encoder), an architecture that deliberately restricts analysis to short windows. Our results demonstrate that this constrained approach matches or exceeds the performance of context-heavy models (XLS-R (Babu et al., 2022)) across diverse languages while using an order of magnitude fewer parameters and providing clean, context-independent phonemic representations suitable for allophone analysis.

## 2 Contextless Universal Phoneme Model

Analysis of our evaluation datasets (Table 2) shows phoneme durations averaging 80ms (range: 62-



Figure 1: The windowing approach restricts the model's context for better localized learning, therefore, generalizing better across languages without learning longer patterns.

107ms), consistent with Crystal and House's findings of 70-120ms for English phonemes (Crystal and House, 1988). Our architecture processes acoustic features through Conv1D layers at 13.1ms per frame, with a 120ms window and 80ms stride to capture 1-2 phonemes per window. This approach provides precise frame-level analysis while maintaining phoneme-level context, departing from traditional methods that rely on broader windows. To preserve acoustic continuity across overlapping windows, we implement a cosine-based weighting mechanism for feature fusion. The complete model architecture is illustrated in Figure 1, with detailed specifications provided in Table 1.

### 2.1 Window Slicer

The Window Slicer module addresses the fundamental challenge of processing continuous speech signals by segmenting raw waveforms ($16\,\text{kHz}$) into overlapping windows. This design enables localized feature extraction while preserving temporal continuity at boundaries. Using a $120\,\text{ms}$ window size with an $80\,\text{ms}$ stride provides sufficient context for phonetic events while reducing computational complexity from $\mathcal{O}(T^2)$ to $\mathcal{O}(W^2)$, where $T$ is the total sequence length and $W$ is the window size.

Given an input audio signal $\mathbf{x} \in \mathbb{R}^{B \times T}$, where $B$ is the batch size and $T$ is the total number of samples in the input sequence ($T = \text{sample\_rate} \times \text{duration}$):

$$w_{b,i}(t) = x_b(t + is), \quad t \in [0, W-1] \quad (1)$$

where $b \in [0, B-1]$ is the batch index, $i \in [0, N-1]$ is the window index, $t$ is the time index within each window, $W = 1920$ is the window

size (120 ms × 16 kHz), $s = 1280$ is the stride length (80 ms × 16 kHz), and $N = \lfloor \frac{T-W}{s} \rfloor + 1$ is the number of windows.

## 2.2 Feature Extractor

Drawing from raw waveform processing techniques (Dai et al., 2017; Schneider et al., 2019), our feature extraction stage implements a hierarchical CNN architecture that processes raw waveforms directly. This design, detailed in Table 1, captures increasingly abstract representations while maintaining computational efficiency. Following the success of Squeeze-and-Excitation Networks (Hu et al., 2018) in speech recognition (Han et al., 2020), we incorporate adaptive channel-wise recalibration through frequency attention. The architecture separates temporal and spectral processing streams, inspired by multi-stream approaches (Han et al., 2021), to capture both evolving acoustic patterns and frequency relationships.

## 2.3 Windowwise Transformer

Our transformer encoder layers process independent fixed windows instead of the whole clip, modifying the contextual processing of standard transformers (Vaswani et al., 2017). This approach represents a departure from traditional speech transformers by restricting context to local windows, ensuring that phoneme recognition decisions rely on relevant local context. Our preliminary experiments showed a tendency to overfit with larger transformer layers, leading us to maintain a light architecture (13M parameters for transformer) with a high dropout of 0.25. For comparison, the XLSR model (Conneau et al., 2021) has over 300M parameters.

## 2.4 Classification and Window Stitching

The final stage of our pipeline consists of classification and temporal integration. The transformer outputs first undergo classification through a two-layer neural network, which maps the high-dimensional representations to phoneme logits. This classifier is designed to untangle complex phonetic representations while maintaining computational efficiency. To ensure temporal coherence across window boundaries, we implement a cosine-based weighting scheme:

$$\tilde{y}(b,t,c) = \frac{\sum_k \cos(\pi t/F_w - \pi/2) \cdot y_k(b,t,c)}{\sum_k \cos(\pi t/F_w - \pi/2) + \epsilon},$$
$$t \in [0, F_w] \tag{2}$$

where $y_k(b,t,c)$ represents the logit from window $k$ for batch $b$, time $t$, and class $c$. This weighted stitching approach enables effective recognition of phonemes shorter than the window length while preserving temporal coherence.

## 3 Experimentation

We experiment with both supervised and self-supervised learning for the proposed model. First, we evaluated our model architecture using labeled speech and phoneme sequences. Then, we adapted the same architecture for self-supervised pretraining using vector quantization projections as targets, following a wav2vec-inspired approach. For baseline comparison, we use the XLS-R (Babu et al., 2022) 300M architecture with an additional linear classification layer. In non-pretrained evaluations, we reset XLS-R's parameters, while for pretrained evaluations, we fine-tune the off-the-shelf model with optional feature extraction layer freezing. The experimental pipeline remains consistent across all tests, varying only in context length (120ms, 160ms, 360ms, or complete words), model selection (XLS-R or CUPE), XLS-R parameter reset status, and feature extraction layer freezing status.

### 3.1 Datasets

We evaluate our model on three diverse speech corpora:
**(1) FLEUR** (Few-shot Learning Evaluation of Universal Representations of Speech) (Conneau et al., 2023): Used exclusively for self-supervised pretraining, comprising 5 hours of audio data from each of 102 languages. Table 2 reports trimmed durations excluding leading and trailing silences.
**(2) Multilingual Spoken Words Corpus (MSW)** (Mazumder et al., 2021): Contains isolated words from Mozilla Common Voice. We use 32 high-resource languages for training (10-hour limit per language) and 6 low-resource languages (lt, mt, ia, sk, ka, as) for evaluation. Twelve languages were excluded due to incompatibility with espeak-NG (esp, 2022), the tool we used to generate IPA phoneme sequences from text.

Table 1: Detailed architecture specifications of the CUPE model with 30M trainable parameters.

| Layer | Output Shape | Parameters | TR | RF | Other Details |
|---|---|---|---|---|---|
| →Window | (B, 1, 1920) | - | 80ms | 120-360ms | Speech waveforms at 16kHz |
| Conv1D-1 | (B, n, 275) | k=15, s=7, p=7 | 13.1ms | 150ms | + BatchNorm + GELU + D(0.1) |
| Conv1D-2 | (B, 2n, 55) | k=11, s=5, p=5 | 1.9ms | 21.3ms | + BatchNorm + GELU + D(0.1) |
| Conv1D-3 | (B, 4n, 19) | k=7, s=3, p=3 | 0.4ms | 4.2ms | + BatchNorm + GELU + D(0.1) |
| Conv1D-4 | (B, 8n, 10) | k=5, s=2, p=2 | 0.1ms | 1.3ms | + BatchNorm + GELU + D(0.1) |
| Freq. Attention | (B, 8n, 10) | k=1, s=1, p=0 | 0.1ms | 0.6ms | ⊙ AvgPool+Conv1D+Sigmoid |
| Temporal Stream (TS) | (B, 8n, 10) (B, 8n, 10) | k=7, s=1, p=3 k=3, s=1, p=1 | 13.1ms | ±11 frames | 2×Conv1d, g=8, BN, GELU |
| Spectral Stream (SS) | (B, 12n, 10) (B, 8n, 10) | k=1, s=1, p=0 k=1, s=1, p=0 | 13.1ms | 1 frame | 2×Conv1d, g=8, BN, GELU |
| Fusion | (B, 8n, 10) | k=1, s=1, p=0 | 13.1ms | | Concat (TS, SS) + 1x1 Conv + BN + GELU |
| Transformer | (B, 10, 512) | $F_w$=10@120ms | 13.1ms | Full window | 4 layers, 8 heads, Pre-norm, D(0.25) |
| ↪FT-Classifier | (B, 10, C) | D=0.25 | 13.1ms | Full window | Supervised only (512→2048→C) |
| ↪PT-Projection | (B, 10, 256) | D=0.25 | 13.1ms | Full window | Unsuperv. only (512→2048→256) |

TR: Temporal Resolution, RF: Receptive Field, B: batch size, k: kernel, s: stride, p: padding, n: base channels (256)

**(3) UCLA Phonetic Corpus (UPC)** (Li et al., 2021b): Features phonetically transcribed speech from 95 languages. We partition this dataset based on language overlap with XLS-R pretraining and FLEUR: UPC-eval contains 64 previously unseen languages, while UPC-seen includes 25 languages present in both pretrained XLS-R and FLEUR. The remaining six languages (fa, ig, kea, ab, eu, haw), exclusive to either XLS-R or FLEUR, serve as validation data during supervised training.

Table 2 summarizes the dataset statistics. The corpora differ significantly in language family distribution and recording conditions. MSWC and FLEUR predominantly feature Indo-European languages by duration, while UPC comprises 48.5% African languages. MSWC offers diverse speakers and recording environments per language, whereas UPC contains just 60 utterances per language, typically from a single speaker in consistent recording conditions.

## 3.2 Pre-Processing

One of the fundamental challenges in creating a universal phoneme recognition system is accommodating unique phoneme inventories across languages. Prior work has explored two main approaches: probabilistic matching (Liu et al., 2023; Li et al., 2021a), which maps phonemes from new languages to acoustically similar training phonemes, and attribute-based decomposition (Glocker et al., 2023), which reconstructs language-specific phonemes from 35 articulatory attributes using the target language's IPA inventory. While both enable automated adaptation to new languages, they face tradeoffs in precision and feature completeness. Our approach instead employs

systematic manual mapping of rare phonemes to standardized phoneme classes, prioritizing perceptual similarity over articulatory phonological relationships. Our mapping preserves high-frequency palatalized consonants ($t^j$, $n^j$, $r^j$) while merging less frequent ones, maintains perceptually distinct vowel contrasts (e.g., ʌ vs ə, ɪ vs i), keeps length distinctions for frequent vowels (aː, eː, iː, oː, uː), and maps rare phonemes to frequent counterparts based on confusion patterns (e.g., ɒ → a, ɕ → k). For affricates, we maintain distinct representations for common ones (ts, tʃ, dʒ) while simplifying rare variants (pf → f), guided by both frequency and confusion patterns. The mapping dictionary is publicly available along with the source code to facilitate adoption and improvement.

Table 2: Datasets' details. $L_n$: total languages or lang code.

| Set | $L_n$ | Hrs | WD(std) | PPW(std) | $U/C$ |
|---|---|---|---|---|---|
| MSWCtrain | 32 | 181 | 0.80(0.12) | 6.30(1.45) | 803/65 |
| MSWCeval | 6 | 15.6 | 0.82(0.12) | 6.39(1.34) | 117/56 |
| Lithuanian | lt | 5.2 | 0.87(0.12) | 6.58(1.34) | 66/42 |
| Maltese | mt | 4.9 | 0.77(0.11) | 6.25(1.30) | 56/38 |
| Interlingua | ia | 3.17 | 0.84(0.12) | 5.98(1.26) | 29/29 |
| Slovak | sk | 1.37 | 0.88(0.11) | 6.77(1.3) | 43/38 |
| Georgian | ka | 0.87 | 0.82(0.11) | 6.97(1.33) | 34/28 |
| Assamese | as | 0.05 | 0.77(0.12) | 5.80(1.21) | 31/26 |
| UPC-eval | 67 | 0.82 | 0.93(0.20) | 5.01(1.53) | 237/59 |
| UPC-seen | 28 | 0.56 | 0.89(0.22) | 4.89(1.32) | 221/60 |
| FLEURS | 102 | 455 | - | - | - |

WD: avg. Word Duration (s), PPW: avg. Phonemes-Per-Word
$U$: Unmapped unique phonemes, $C$: Mapped phoneme classes.

## 3.3 Supervised Training

For each window, the model generates frame-level logits (10 frames per 120ms window, 28 frames for 360ms), which are stitched into continuous phoneme sequences. Training uses CTC

loss (Graves, 2012) with an additional silence-awareness term:

$$\mathcal{L}1 = \mathcal{L}\mathrm{ctc} + \alpha_s \mathcal{L}_{\mathrm{sil}} \qquad (3)$$

$$\mathcal{L}_{\mathrm{sil}} = \frac{1}{B} \sum_{t,b} (0.5\tilde{y}_b^t M_s^t + 0.1\tilde{y}_b^t(1 - M_s^t)) \quad (4)$$

where $\tilde{y}_b^t$ is the blank token probability, $M_s^t$ is the silence mask, $B$ is batch size, and $\alpha_s$ (default 0.01) balances silence detection with phoneme recognition.

Our training pipeline optimizes for efficient learning through several mechanisms: AdamW optimizer with OneCycleLR scheduling, gradient norm clipping at threshold $\tau = 1.0$, and mixed-precision BF16 training for balanced efficiency and numerical stability. We trained all models on MSWC-train using a batch size (B) of 300 words until validation PER showed no further improvement, requiring 20 epochs and approximately 7 hours on two A6000 GPUs. The trained models and source code are available online[0], with results presented in Table 3.

### 3.4 Self-supervised Pre-Training

For self-supervised pre-training, we modify CUPE by replacing the FT-Classifier with a prediction head (two projection layers with residual connections, layer normalization, GELU activation, and dropout 0.1) while being projected to a 256-dimensional feature space. The core architecture remains unchanged.

The pre-training uses masked prediction on 120ms windows (80ms stride), masking 40% of features based on energy profiles and acoustic boundaries, with per-batch constraints of 10-80%. A vector quantizer with 256-entry codebook serves as training target, using EMA updates (decay 0.99) and Laplace smoothing. The training objective combines reconstruction loss (smooth L1), contrastive loss with curriculum learning, codebook diversity loss, and similarity regularization.

Optimization uses AdamW (weight decay 0.05) with hierarchical learning rates (encoder: 5e-4, quantizer: 1e-3, prediction head: 1.5e-3) and one-cycle scheduling (15% warmup, momentum 0.8-0.9). For evaluation, we freeze the feature extractor, replace the prediction head with classification layers, and fine-tune only the transformer and FT-Classifier components. We similarly evaluate XLS-R with both full and frozen-backbone fine-tuning.

### 3.5 Results

#### 3.5.1 Evaluation Metrics

We decoded model outputs using Greedy Best-First Search and evaluated using Phoneme Error Rate (PER), Ground-truth Probability (**GP**), and F1-score. GP and F1 are computed after optimal alignment of true and predicted sequences, excluding insertions and deletions. While PER assigns a full penalty (+1) for any substitution, insertion, or deletion, it doesn't measure the near-misses. We introduce GP (**GPm** for macro, **GPw** for class-weighted) to better evaluate fine-grained phonemic distinctions like duration variants (i/iː) and vowel contrasts (æ/a) that are preserved in our approach rather than merged. GP measures the model's probability assignment to ground-truth classes at aligned time steps. It can be intuitively understood as the proximity to truth, or conversely, the inverse of the distance from truth. This proximity measure instead of PER is more important for judging the quality of embeddings for latent tasks.

Detailed analysis of model behavior is provided in Appendix A. The confusion matrix in Figure 2 shows that contextless recognition errors follow phonetically meaningful patterns, with confusions primarily occurring between acoustically similar sounds (e.g., front vowels, voiced/voiceless consonant pairs) rather than random misclassifications. The phoneme probability distributions over time (Figure 3) illustrate CUPE's temporal resolution capabilities, showing distinct probability peaks corresponding to ground truth phonemes and smooth transitions between adjacent sounds.

### 3.6 Key Insights and Limitations

Looking at Table 3, CUPE demonstrates remarkable cross-lingual generalization despite having a fraction of XLSR's parameters. While the 360ms model shows slightly better PER, this can be misleading due to class imbalances - it performs better on long and common vowels like /aː/ but struggles with short but rare phonemes, highlighting why GPm is a more balanced metric. Note that both 360ms and 120ms models have the same frame length of 16ms, the only difference is the context length. The significant performance difference in

---

[0]https://github.com/tabahi/contexless-phonemes-CUPE

Table 3: Evaluation metrics (%) for two architectures, XLSR (300M) & CUPE (30M), trained on MSWC-train without pretraining.

| Model:Context | Evaluation on MSWC-eval | | | | Zero-shot PER on individual langs | | | | | | Zero-shot evaluation on UPC-eval | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | PER↓ | GPm | GPw | F1 | lt | mt | ia | sk | ka | as | PER↓ | GPm | GPw | F1 |
| XLSR:word | 49.9 | 35 | 51.7 | 60.6 | 59.5 | 48.7 | 37 | 45.3 | 48.4 | 65.8 | 66.5 | 31.2 | 51.7 | 52.9 |
| XLSR:120ms | 52.6 | 34 | 52.1 | 59.9 | 61.1 | 49.9 | 42.9 | 52.3 | 50.4 | 63.7 | 66.3 | 31.6 | 51.1 | 54.9 |
| CUPE:word | 46.4 | 39 | 55.1 | 63 | 54.5 | 47.1 | 33.1 | 42.5 | 44 | 60.5 | 58.8 | 32.9 | 52.5 | 58.3 |
| CUPE:360ms | **44.8** | 38.3 | 56.5 | 62.6 | **53.8** | 45.2 | **30.8** | **39.7** | 42.5 | 60.9 | **52.2** | 34.7 | 53.1 | 61 |
| CUPE:160ms | 47.8 | 36 | 55 | **64.8** | 57.2 | 46.2 | 36.2 | 45.2 | 44 | 60.7 | 57.5 | 32.9 | 54.1 | 58.8 |
| CUPE:120ms | 45.9 | **40** | **57.5** | 64.5 | 54.6 | **45** | 33.9 | 43.6 | **42.2** | **60.2** | 56.9 | **35.1** | **56.4** | **67.7** |

Table 4: Evaluation metrics (%) for pre-trained models CUPE-PT (30M, pretrained on FLEURS), fine-tuned on MSWC-train, compared with XLSR (300M, off-the-shelf pretrained on 128 languages) with or without frozen backbone (FB) feature extractor. The top 4 rows show the results for contextless (120ms) models, the bottom 4 rows show results for word-context models for reference. Only the UPC-eval languages are unseen languages for zero-shot evaluation.

| Model:Context | Eval. on MSWC-eval | | | | PER↓ on individual langs (seen) | | | | | | UPC-eval | | UPC-seen | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | PER↓ | GPm | GPw | F1 | lt | mt | ia | sk | ka | as | PER↓ | GPm | PER↓ | GPm |
| *120ms Context Models* | | | | | | | | | | | | | | |
| FB-XLSR | 65.8 | 36.2 | **60.2** | 51.4 | 69.6 | 70.7 | 55.6 | 55.0 | 63.6 | 84.5 | 66.3 | **43.5** | 67.8 | 43.5 |
| FB-CUPE-PT | 49.8 | 34.9 | 53.0 | 60.5 | 59.3 | 47.5 | 38.5 | 48.5 | 44.6 | 61.4 | 66.5 | 35.4 | 69.7 | 38.2 |
| XLSR | 52.2 | 38.2 | 56.7 | 62.3 | 60.9 | 48.5 | 43.0 | 51.8 | 50.8 | 67.1 | 63.6 | 37.8 | 60.9 | 45.8 |
| CUPE-PT | **45.6** | **41.2** | 58.1 | **64.0** | 54.5 | 45.2 | 33.5 | 47.9 | 43.6 | 62.1 | **56.2** | 36.4 | **57.6** | **44.2** |
| *Word Context Models* | | | | | | | | | | | | | | |
| FB-XLSR | **43.5** | **40.0** | 58.4 | 68.1 | 53.9 | 42.8 | 30.1 | 38.2 | 37.3 | 55.3 | 66.9 | **48.5** | 70.3 | 43.4 |
| FB-CUPEPT | 70.4 | 1.9 | 29.6 | 54.3 | 73.6 | 65.1 | 71.2 | 77.1 | 70.5 | 62.5 | 69.0 | 3.2 | 73.2 | 2.7 |
| XLSR | 46.6 | 36.3 | 53.6 | 66.7 | 56.7 | 44.6 | 35.5 | 39.7 | 44.2 | 63.8 | **46.9** | 39.8 | **46.0** | **46.4** |
| CUPE-PT | 46.1 | 38.1 | 56.1 | 61.4 | 54.2 | 45.6 | 35.5 | 41.7 | 42.6 | 60.4 | 56.8 | 37.9 | 54.0 | 46.2 |

UPC evaluations, even when XLSR:120ms uses the same windowing pipeline, suggests that model's heavy size could be an overfitting liability.

Table 4 reveals that while XLSR with a frozen feature extractor achieves better overall metrics, CUPE maintains competitive performance under significant constraints. Notable observations include XLSR's degraded performance on UPC with frozen features and CUPE's sharp performance drop with word-context windows, perhaps due to having to learn more phonemes per window while most parameters are frozen. The completely unfrozen CUPE model's results mirror those in Table 3 even though the learning rate was set 10 times less for fine-tuning. The best contextless model, CUPE-PT:120ms, does not perform as well as pre-trained XLSR with full word context, indicating that additional context and parameters benefit large-scale pretraining. Nevertheless, CUPE's effectiveness with frozen feature extractors shows that essential phonetic information is learned by the feature extractor within brief temporal windows during pretraining. Another sharp degradation is noticeable for CUPE-PT word context compared to 120ms; it is possibly due to 30M parameters being not enough for longer sequences (1000ms vs 120ms).

Our approach faces several limitations in its current form. The fixed 120 ms window presents inherent trade-offs in phoneme recognition: too long for short stop consonants and insufficient for capturing long phonemes fully. The model shows the best recall of stop consonants, but the worst recall of infrequent vowels. This issue is particularly evident in languages with contrastive length distinctions, where the model struggles to maintain consistent performance across different phoneme durations.

The performance gap between supervised and pre-trained+fine-tuned results points to architectural limitations in both the projection mechanism and loss objectives. The current projection approach may not optimally preserve phonetic features during self-supervised learning, while the loss objectives could better reflect the hierarchical nature of phonemic contrasts. Additionally, the relatively modest size of the model (30M parameters) may limit its capacity to capture the full complexity of cross-linguistic phonetic variations. Additionally, our systematic mapping of rare phonemes, while practical, may obscure certain phonological contrasts. Although we achieve competitive results on the UCLA Phonetic Corpus, direct comparisons with methods such as Epitran (Li et al., 2021a) and

Table 5: Zero-shot PER comparison on UPC (UCLA Phonetic Corpus) with other works. Our CUPE:120ms results are fine-tuned on language splits matched to each baseline study for fair comparison, which differ from the UPC-eval/UPC-seen partitions in Tables 3-4. Direct performance comparison is limited due to different phoneme mapping systems. $L_n$ = number of unseen test languages (of 95).

| Study | $L_n$ | ↓ PER (%) | Phoneme Inventory Approach |
|---|---|---|---|
| (Li et al., 2021a) | 47 | 51.2 | Epitran+Allovera+Panphon |
| Ours | 47 | 46.1 | Systematic mapping to 65 classes |
| (Liu et al., 2023) | 10 | 64.7 | Direct UPC inventory |
| Ours | 10 | 44.1 | Systematic mapping to 65 classes |
| (Li et al., 2022) | 77 | 64.2 | Bayesian tree-based estimation |
| Ours | 77 | 48.6 | Systematic mapping to 65 classes |
| (Glocker et al., 2023) | 84 | 45.62 | 35 articulatory attribute system |
| Ours | 84 | 48.98 | Systematic mapping to 65 classes |

Allophant (Glocker et al., 2023) are challenging due to fundamentally different phoneme inventory approaches.

While CUPE demonstrates strong performance in contextless phoneme recognition, several limitations warrant discussion. The model's varying performance across language families suggests potential biases in the feature extraction process that merit further investigation. Some languages with distinct phonological structures or phoneme inventories may require specialized preprocessing or architectural adaptations to achieve optimal performance. Additionally, the fixed 120ms window size, while effective across our evaluation datasets, may not be optimal for all languages or phonetic contexts—some phonemes naturally require longer or shorter temporal windows for accurate characterization.

Most importantly, this work establishes the foundation for more complex speech analysis systems. We have demonstrated how to extract clean embeddings for individual allophones—the next critical step is implementing a sentence-level speech style encoder that learns from these contextless allophone embeddings. Such a system would enable comprehensive analysis of speaker characteristics, accent patterns, and speaking styles while maintaining the interpretability and cross-linguistic generalizability that contextless representations provide.

While our approach achieves competitive results on the UCLA Phonetic Corpus compared to existing methods listed in Table 5, these comparisons should be interpreted cautiously - each method uses fundamentally different phoneme inventory systems, from Epitran's probabilistic mappings (Li et al., 2021a) to Allophant's 35 articulatory attributes (Glocker et al., 2023), making direct performance comparisons less meaningful. Our choice of 65 systematically mapped classes represents a

different trade-off between granularity and generalization. The 65 class system is pragmatic implementation which can be expanded depending on the dataset. We selected 65 phonemes by empirically analyzing their occurrence across MSWC's 50 languages, including only those that appeared at least 10,000 times. While phoneme mapping can further reduce the number of classes, our findings show that the impact on error rate is limited. For instance, when we applied broad phoneme group mapping to reduce the set to just 15 phonemes, the PER on MSWC-eval dropped from 0.45 to 0.40.

## 4 Conclusion

Through this work, we have demonstrated that effective universal phoneme recognition can be achieved using brief 120ms windows of speech input. Our CUPE model achieves competitive performance while requiring an order of magnitude fewer parameters than current approaches. The model's success in cross-lingual generalization validates our core finding that essential phonetic information can be captured through focused analysis of brief speech segments. These results provide compelling evidence that extensive temporal context is not a requirement for robust speech processing tasks. While our approach has some limitations, particularly with very long phonemes and limited phoneme inventory, it opens promising directions for lightweight, language-agnostic speech processing systems. CUPE's effectiveness has significant implications for real-world applications, from low-latency speech recognition and ASR self-learning to speech pathology diagnostics. Our results indicate that future speech processing systems may benefit from focusing on fundamental acoustic patterns rather than extensive contextual dependencies.

# References

2022. eSpeak NG (version 1.50). https://github.com/espeak-ng/espeak-ng/tree/1.50. [Online].

A. Babu, C. Wang, A. Tjandra, K. Lakhotia, Q. Xu, N. Goyal, K. Singh, P. von Platen, Y. Saraf, J. Pino, A. Baevski, A. Conneau, and M. Auli. 2022. XLS-R: Self-supervised Cross-lingual Speech Representation Learning at Scale. In *Proc. Interspeech 2022*, pages 2278–2282.

Alexei Baevski, Yuhao Zhou, Abdelrahman Mohamed, and Michael Auli. 2020. wav2vec 2.0: A framework for self-supervised learning of speech representations. *Advances in neural information processing systems*, 33:12449–12460.

Alexis Conneau, Alexei Baevski, Ronan Collobert, Abdelrahman Mohamed, and Michael Auli. 2021. Unsupervised cross-lingual representation learning for speech recognition. *Interspeech 2021*.

Alexis Conneau, Min Ma, Simran Khanuja, Yu Zhang, Vera Axelrod, Siddharth Dalmia, Jason Riesa, Clara Rivera, and Ankur Bapna. 2023. Fleurs: Few-shot learning evaluation of universal representations of speech. In *2022 IEEE Spoken Language Technology Workshop (SLT)*, pages 798–805. IEEE.

Thomas H Crystal and Arthur S House. 1988. Segmental durations in connected-speech signals: Current results. *The journal of the acoustical society of America*, 83(4):1553–1573.

Wei Dai, Chia Dai, Shuhui Qu, Juncheng Li, and Samarjit Das. 2017. Very deep convolutional neural networks for raw waveforms. In *2017 IEEE international conference on acoustics, speech and signal processing (ICASSP)*, pages 421–425. IEEE.

Katharina Glocker, Ans Herygers, and Michael Georges. 2023. Allophant: Cross-lingual phoneme recognition with articulatory attributes. In *Proceedings of Interspeech 2023*, pages 2258–2262.

Alex Graves. 2012. Connectionist temporal classification. In *Supervised sequence labelling with recurrent neural networks*, pages 61–93. Springer.

Kyu J Han, Jing Pan, Venkata Krishna Naveen Tadala, Tao Ma, and Dan Povey. 2021. Multistream cnn for robust acoustic modeling. In *ICASSP 2021-2021 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, pages 6873–6877. IEEE.

Wei Han, Zhengdong Zhang, Yu Zhang, Jiahui Yu, Chung-Cheng Chiu, James Qin, Anmol Gulati, Ruoming Pang, and Yonghui Wu. 2020. ContextNet: Improving Convolutional Neural Networks for Automatic Speech Recognition with Global Context. In *Proc. Interspeech 2020*, pages 3610–3614.

Jie Hu, Li Shen, and Gang Sun. 2018. Squeeze-and-excitation networks. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 7132–7141.

Xinjian Li, Juncheng Li, Florian Metze, and Alan W Black. 2021a. Hierarchical phone recognition with compositional phonetics. In *Interspeech*, pages 2461–2465.

Xinjian Li, Florian Metze, David R Mortensen, Alan W Black, and Shinji Watanabe. 2022. Phone inventories and recognition for every language. In *Proceedings of the Thirteenth Language Resources and Evaluation Conference*, pages 1061–1067.

Xinjian Li, David R Mortensen, Florian Metze, and Alan W Black. 2021b. Multilingual phonetic dataset for low resource speech recognition. In *ICASSP 2021-2021 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, pages 6958–6962. IEEE.

Qianying Liu, Zhuo Gong, Zhengdong Yang, Yuhang Yang, Sheng Li, Chenchen Ding, Nobuaki Minematsu, Hao Huang, Fei Cheng, Chenhui Chu, et al. 2023. Hierarchical softmax for end-to-end low-resource multilingual speech recognition. In *ICASSP 2023-2023 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, pages 1–5. IEEE.

Mark Mazumder, Sharad Chitlangia, Colby Banbury, Yiping Kang, Juan Manuel Ciro, Keith Achorn, Daniel Galvez, Mark Sabini, Peter Mattson, David Kanter, et al. 2021. Multilingual spoken words corpus. In *Thirty-fifth Conference on Neural Information Processing Systems Datasets and Benchmarks Track (Round 2)*.

Steffen Schneider, Alexei Baevski, Ronan Collobert, and Michael Auli. 2019. wav2vec: Unsupervised Pre-Training for Speech Recognition. In *Proc. Interspeech 2019*, pages 3465–3469.

Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N. Gomez, Lukasz Kaiser, and Illia Polosukhin. 2017. Attention is all you need. *Proceedings of the 31st International Conference on Neural Information Processing Systems (NIPS 2017)*.

# A Confusion Heatmaps



Figure 2: Confusion matrix for contextless phoneme recognition on MSWC-eval dataset using CUPE:120ms model trained on MSWC-train. The heatmap shows predicted phonemes (x-axis) versus ground truth phonemes (y-axis), with color intensity indicating count frequency on a logarithmic scale. The 'Un' counts show the unaligned trues or predictions (i.e., the true sequence had a phoneme that didn't exist or aligned in the predicted sequence and vice-versa). The matrix reveals systematic confusion patterns, with darker cells along the diagonal indicating correct predictions. Notable off-diagonal clusters highlight acoustically similar phoneme pairs that are challenging for contextless recognition, such as front vowels, central vowels, and voiced/voiceless consonant pairs. The sparse structure demonstrates that most confusions occur within phonetically related categories rather than across distant phoneme classes.

Figure 3: Phoneme probability distributions over time for an example utterance using CUPE:120ms model. The top panel shows a heatmap of phoneme probabilities (y-axis) across time frames (x-axis), with color intensity representing probability values. Ground truth phoneme alignments are displayed at the bottom with text. The visualization demonstrates the model's ability to capture temporal phoneme transitions in contextless recognition, with clear probability peaks corresponding to ground truth phonemes. Notable patterns include smooth transitions between phonemes within words and distinct silence regions (SIL) between words, highlighting the model's temporal resolution at 13ms.

# The Impact of Annotator Personas on LLM Behavior Across the Perspectivism Spectrum

**Olufunke O. Sarumi[1], Charles Welch[2], Daniel Braun[1], Jörg Schlötterer[1,3]**

[1]University of Marburg, [2]McMaster University, [3]University of Mannheim

{sarumio,daniel.braun,joerg.schloetterer}@uni-marburg.de[1], cwelch@mcmaster.ca[2]

## Abstract

In this work, we explore the capability of Large Language Models (LLMs) to annotate hate speech and abusiveness while considering pre-defined annotator personas within the strong-to-weak data perspectivism spectra. We evaluated LLM-generated annotations against existing annotator modeling techniques for perspective modeling. Our findings show that LLMs selectively use demographic attributes from the personas. We identified prototypical annotators, with persona features that show varying degrees of alignment with the original human annotators. Within the data perspectivism paradigm, annotator modeling techniques that do not explicitly rely on annotator information performed better under weak data perspectivism compared to both strong data perspectivism and human annotations, suggesting LLM-generated views tend towards aggregation despite subjective prompting. However, for more personalized datasets tailored to strong perspectivism, the performance of LLM annotator modeling approached, but did not exceed, human annotators.

## 1 Introduction

Perspectivism in Natural Language Processing (NLP) aims to preserve the spectrum of opinions held by annotators in corpora (Cabitza et al., 2023). Dataset annotation for this purpose often uses a descriptive paradigm (Rottger et al., 2022), involving minimal instructions and multiple annotators providing labels for every corpus sentence to capture diverse viewpoints. The number of annotators involved can range significantly, from a minimum of 2 to 2500 or more (Plepi et al., 2022; Frenda et al., 2024).

Most traditional approaches aggregate labels to obtain a single majority label (Davani et al., 2022; Aroyo and Welty, 2015), which is commonly used for training models. However, the perspectivist approach argues that critical information is lost when labels are aggregated. More importantly, the opinions of the minority, which may represent a significant population, are undermined, leading to under-representation and overshadowing of nuances inherent in the dataset. This is crucial because people's views and opinions are indeed shaped by different socio-cultural, demographic, economic, and experiential backgrounds (Akhtar et al., 2021; Almanea and Poesio, 2022; Demszky et al., 2020; Kennedy et al., 2022). These factors impact how individuals perceive, interpret, and respond to various topics, making it unrealistic to assume everyone shares similar views on the same subject. Recognizing and reflecting opinion differences in our models is therefore important for developing socially aware NLP systems, treating disagreements not as errors but as distinct perspectives. To address this, models have been developed that can learn from such disaggregated labels (Leonardelli et al., 2023; Sullivan et al., 2023; Vitsakis et al., 2023; García-Díaz et al., 2023; Cui, 2023; Xu et al., 2024).

Furthermore, while some disagreements stem from different perspectives, other factors also cause disagreement in data annotations, including temporal factors, annotator inconsistencies, uncertainty, ambiguities, lack of task understanding, or a perfunctory approach to annotation (Fleisig et al., 2024). When modeling perspectives obtained from subjective tasks, these perspectives are often mixed with noise and errors, raising the question of whether true perspectives or merely annotator inconsistencies have been modeled. Some literatures have quantified these uncertainties to a minimal extent (Klemen and Robnik-Šikonja, 2022; Davani et al., 2022).

In this work, we aimed to investigate how existing annotator modeling techniques would behave when trained on deterministic LLM-generated annotations, in contrast to earlier work that explored modeling individual human annotators' perspectives using disaggregated labels. We generated

new annotations for the HS-Brexit and ConvAbuse datasets using Llama2-13B, guided by persona-based prompting derived from annotator information provided by the original authors.

In generating these annotations, we implemented two perspectivism approaches: *strong* and *weak* data perspectivism. Weak perspectivism, also known as reduced perspective, involves considering multiple labels which are ultimately aggregated into one, representing a group opinion. Strong perspectivism, by contrast, utilizes and retains all distinct labels from training through evaluation (Cabitza et al., 2023; Frenda et al., 2024).

Our findings show that LLMs struggle to generate responses as diverse as humans, even with diverse personas. They still partially align with human annotations but tend to pick up only selected persona features. Furthermore, we identified latent annotation prototypes shared by multiple human annotators. These alignment patterns vary across datasets and perspectivism strategies: for instance, HS-Brexit with contrasting demographic attributes shows stronger alignment with human annotations under weak perspectivism, whereas ConvAbuse demonstrates closer alignment with human annotations when strong data perspectivism is used, involving highly personalized and overlapping persona features.

## 2 Related Work

The first part of this section addresses how Large Language Models (LLMs) have been used to generate different perspectives and their ability to adopt an assigned persona. It also highlights the lack of connection between perspectivism, based on defined personas and annotations in subjective tasks. The second part focuses on the use of LLMs as annotators, examining their ability to generate discrete multiple labels, identifying the lack of persona-based labeling, and replicating human annotation behavior to enable alignment with human annotations.

### 2.1 LLMs in Perspectivism and Adopting Personas

LLMs have been explored for their ability to simulate diverse human perspectives. Subjective tasks often involve annotators with different backgrounds, leading to divergent opinions which often reflect demographic variation, different and substantial opinions, these make label aggregation in-

adequate (Rottger et al., 2022). Some works argue that LLMs naturally contain persona traits, as they are trained in large corpora, often culled from social networks that contain crowd-sourced data rich with diverse viewpoints (Hu and Collier, 2024; Vitsakis et al., 2023). For example, Hayati et al. (2024) showed that it is possible to generate multiple perspectives from LLMs and quantify the maximum number of perspectives derivable from an LLM. However, the influence of persona prompting remains debated and the influence of specific persona traits remains underexplored (Beck et al., 2024; Sun et al., 2025). Hu and Collier (2024) suggests that personas have minimal effect on LLM outputs, whereas a psycholinguistic research found that LLMs can generate human-like outputs, even surpassing humans in turing experiments, yet exhibit unnaturally high accuracy that is not possible within human populations (Aher et al., 2023). Furthermore, Wang et al. (2024) found that LLMs risk homogenizing or misrepresenting marginalized identity groups, particularly when asked to simulate them. These challenges highlight the difficulty in separating the LLM's inherent persona from externally applied persona prompts. Despite this, prompting LLMs with well-defined personas, particularly those grounded in demographic traits from existing datasets, offers a practical way to examine how perspective alignment occurs between machines and humans. However, small variations in prompt configurations can lead to large differences in output, complicating reproducibility and fairness evaluations.

### 2.2 LLM Annotations and Label Generation

Beyond simulating perspectives, LLMs are being explored as direct substitutes for human annotators (Ivey et al., 2024; Bavaresco et al., 2024), especially in settings where collecting human annotations is expensive or slow (Huang et al., 2023; Gligorić et al., 2024). Recent studies have examined the ability of LLMs to generate discrete labels for classification tasks, often using crowd-sourced datasets as benchmarks (Pavlovic and Poesio, 2024a; Gilardi et al., 2023). Gilardi et al. (2023) found that LLMs outperformed crowd-sourced workers in certain annotation tasks, while Pavlovic and Poesio (2024b) demonstrated that adjusting temperature values can control LLM behavior to better simulate annotation disagreement or consistency. These findings suggest that LLMs can be tuned to exhibit behavior similar to individual or aggregated hu-

man annotators. LLMs have also been deployed in replicating prior annotation experiments. For example, Pavlovic and Poesio (2024a) replicated a Learning With Disagreement task (Leonardelli et al., 2023) using GPT-3 but did not incorporate the demographic background of annotators, limiting their insight into perspective-specific agreement. While many experiments rely on LLMs generating explanations or engaging in dialogue-based tasks, fewer works have explored their ability to produce discrete, disaggregated annotations comparable to crowdsourced annotators. Likewise, existing annotator modeling techniques are yet to be fully evaluated on annotations generated by LLMs. The impact of LLM annotations and predefined personas on existing annotator modeling approaches remains unexplored and is a key area we address in our study.

## 3 Dataset

We used two datasets from the SemEval-2023 task on learning with disagreements (Leonardelli et al., 2023) and used Llama2-13B to generate annotations for weak and strong data perspectivism variants resulting in six (6) datasets. Strong perspectivism used prompts tailored to individual persona descriptions, while weak perspectivism used group descriptions to simulate aggregated viewpoints; however, the persona descriptions in each variant were limited to the demographic information and features provided in the original work. All datasets use binary labels for classification. Original dataset statistics are presented in Table 1.

**HS-Brexit** The Hate Speech Brexit (HS-Brexit) dataset (Akhtar et al., 2021) comprises 1,120 tweets concerning Brexit and immigration, annotated for hate speech, aggressiveness, and offensiveness. This dataset features annotations from two distinct groups of three individuals: a target group of Muslims and first- or second-generation immigrants to the UK (also classified as migrants in the original study) and a control group of researchers with a Western background making six annotators in all.

**ConvAbuse** The Conversational Abuse (ConvAbuse) dataset, as described by Cercas Curry et al. (2021), comprises roughly 4,000 English dialogues between users and two conversational agents. These user conversations were labeled by a minimum of three gender studies experts, using a hierarchical annotation system that included categories for presence, severity, and directness of abuse. We

binarized the annotations into two classes, 0 and 1. The ConvAbuse dataset is characterised by eight (8) annotators, each providing a significant number of annotations. Also, not all the 8 annotators labeled every instance contrary to the HS Brexit, but each annotator has annotations.

## 4 Methodology

Firstly, we explore the ability of Llama2-13B to generate discrete binary annotations on the datasets, using defined personas. Secondly, we modeled these personas with existing annotator modeling techniques.

### 4.1 Annotation Generation

For the strong perspectivism variant of the datasets, we prompt Llama2-13B with each text in the original corpus. We extended the dataset with the generated annotations for each corresponding persona, maintaining the original structure of the dataset from the SemEval-2023 task. The *strong* variant uses specific individual descriptions for each persona as seen in Figure 1. In the original ConvAbuse dataset, not all annotators annotated all instances, but in the LLM version, all eight annotators were represented in all instances. We generate annotations at temperatures: 0, 0.1, 0.2, 0.5 and 0.8, for each perspectives. We used the demographic description presented in the original work as guide for our persona features. In *weak* perspectivism, we followed the same approach. Figures 3 and 4 show persona descriptions and Table 2 shows a sample of the prompt used. The prompt and personas are fully described in the Appendices A and B, respectively. Also in Table 3, we show a summary of the data statistics and the variance observed in the inter-annotator agreement K-$\alpha$ as temperature increases.

### 4.2 Annotator Modeling

We trained existing annotator models (Oluyemi et al., 2024; Davani et al., 2022) using the LLM-generated labels, following a classification pipeline originally used with the human-annotated corpus. We replicated these annotator modeling techniques—User Token, Composite Embedding, Composite+User Token Embedding, and Multi-task to model perspectives by modeling annotators, and we also added a text-only implementation without annotator information with SBERT. These techniques used annotator IDs and text, with 6 annotations per instance for the HSBrexit and 8 annotations per

|            | #A | #I    | N                | A/I          | K-$\alpha$ |
|------------|----|-------|------------------|--------------|------------|
| HS-Brexit  | 6  | 1,120 | 1,120.00± 0.00   | 6.00± 0.00   | 0.35       |
| ConvAbuse  | 8  | 4,050 | 1,521.00± 206.91 | 3.00± 0.88   | 0.65       |

Table 1: Original Dataset Statistics by Human Annotators. #A: number of annotators, #I: number of total instances, N: number of annotations per annotator, A/I: annotations per instance, K-$\alpha$: Krippendorff's alpha agreement.



Figure 1: Distribution of Annotations for ConvAbuse and HS-Brexit (Strong Perspective) by Persona

*annotate this statement regarding Brexit. as a group of (persona) in the UK or Europe, you must assign a label to the statement Use '0' if the statement does not express hate towards Brexit, and '1' if it does express hate towards Brexit. Provide only the number as your response without any additional text or explanation for example "####Annotator:"####Annotator:0" or "####Annotator:1"*

Table 2: An example prompt for weak perspectivism in HS-Brexit

instance for the ConvAbuse to capture the *persona* perspectives using the labels obtained from the generations at all temperatures but used the best scores (generally between temperatures 0 and 0.1) in our results and analysis. The annotator ID represents each unique annotator(persona), encoded as embeddings. Each technique uses a different method to generate encodings used to uniquely model the personas. Finally, we compared the performance of these annotator modeling techniques on LLM-generated annotations and human annotations.

## 5 Results

Table 4 presents the F1 scores for the annotator modeling techniques evaluated on both the original and the synthetic datasets. Our analyses show some trends in the performance of these models. In existing results that used human annotations, overall performance was observed on the ConvAbuse dataset. The inter-annotator agreement measured by Krippendorff's alpha was high for ConvAbuse and comparatively lower for the HSBrexit dataset. Interest-

ingly, the Llama2-annotated versions showed significantly higher agreement levels than the original human annotations across all temperature settings, including at a high randomness level (Temperature = 0.8) as seen in Tables 1 and 3. Prior research established that the effectiveness of annotator modeling techniques is largely dependent on the degree of agreement and the number of annotations per annotator (Oluyemi et al., 2024). Specifically, the User-Token modeling approach performs best for datasets with low agreement, while the Composite Embedding + User Token method is optimal for datasets with high agreement. Both methods rely on an explicit naming system, using annotator IDs to individually predict the label outputs for each annotator. However, our results indicate that models without explicit annotator information outperformed others on the Llama2 persona-based datasets. For instance, SBERT, with no annotator information and Composite Embedding- an approach that did not use explicit naming convention (annotator ID) for modeling, both outperformed the best-performing models on HSBrexit and achieved comparable results on ConvAbuse. This suggests that the optimal annotator modeling techniques for human annotations may not be directly transferable or equally effective for data annotated through LLM personas.

| | #A | #I | N | A/I | K-$\alpha$ (Strong) | K-$\alpha$ (Weak) |
|---|---|---|---|---|---|---|
| HS-Brexit | 6 | 1,120 | 1,120.00$\pm$ 0.00 | 6.00$\pm$ 0.00 | 0.58 – 0.81 (T=0.8 – 0) | 0.55 – 0.75 (T=0.8 – 0) |
| ConvAbuse | 8 | 4,050 | 4,050.00$\pm$ 0.00 | 8.00$\pm$ 0.00 | 0.60 – 0.91 (T=0.8 – 0) | 0.62 – 0.93 (T=0.8 – 0) |

Table 3: LLAMA2 Dataset Statistics. #A: number of annotators, #I: number of total instances, N: number of annotations per annotator, A/I: annotati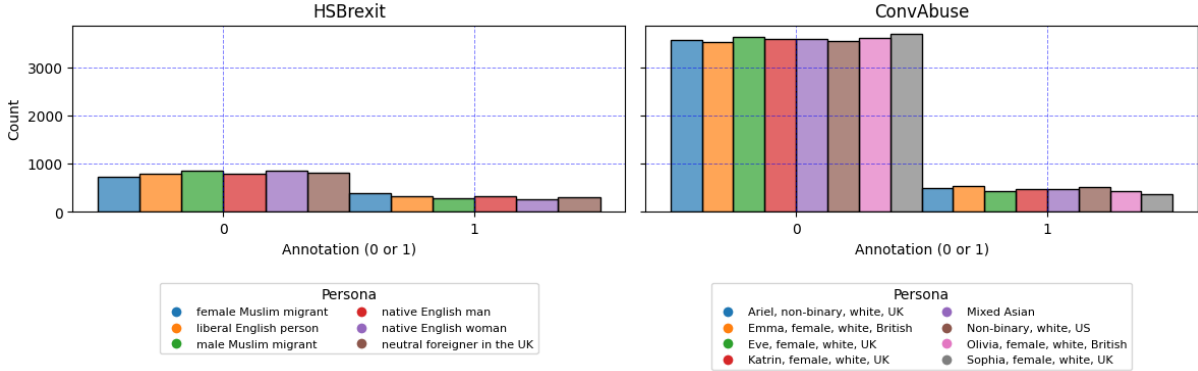ons per instance, K-$\alpha$: Krippendorff's alpha agreement (T=temperature range). The K-$\alpha$ values are presented as a range from temperature 0.8 to 0, that is agreement decreases as temperature increases.

| Method | SBERT | User Token | Composite Embedding | Composite Embedding + User Token | Multi-Tasking |
|---|---|---|---|---|---|
| **Human-annotations** | | | | | |
| HS-Brexit | 68.6 | 77.6 | 67.6 | 77.3 | 71.7 |
| ConvAbuse | 85.9 | 88.5 | 85.8 | 88.6 | 82.3 |
| **LLAMA2-13B strong perspectivism** | | | | | |
| HS-Brexit | **72.2** | 69.4 | **71.8** | 71.2 | 65.1 |
| ConvAbuse | 85.7 | 84.4 | 84.6 | 84.4 | 81.1 |
| **LLAMA2-13B weak perspectivism** | | | | | |
| HS-Brexit | **73.2** | 72.2 | **72.4** | 71.7 | 62.0 |
| ConvAbuse | 85.2 | 83.7 | 83.7 | 81.8 | 79.8 |

Table 4: Model performance based on individual annotator and persona F1 scores. Results for human annotations was adapted from Oluyemi et al. (2024). We reported the best LLM results for temperatures 0 and 0.1.

## 5.1 Strong vs Weak Data Perspectivism in Annotator Modeling

As presented in Tables 6 and 7 of Appendix C, we adapted the two versions of data perspectivism described by (Cabitza et al., 2023) and evaluated the annotator modeling techniques on the datasets. The strong perspectivist approach, which used fine-grained persona profiles, generally produced higher performance that was more aligned with the results from human modeling for the ConvAbuse dataset at temperature 0.1. The weak perspectivism approach, characterized by contrasting group descriptions, showed improved performance over the human version in the HS-Brexit dataset across both strong and weak variants, with a greater improvement observed in the weak, group-based variant. However, this performance increase was exclusively observed in the Composite Embedding and SBERT models without explicit annotator information.

## 5.2 Annotation Quality and Uncertainty

We analyzed the quality of annotations generated by Llama2-13B across a spectrum of temperature parameters. Even at high randomness with temperature set to 0.8, inter-annotator agreement remained high cf. Table 3. The distribution of labels diverged significantly from that of the human annotators. To illustrate this, we compared the label distributions using Probability Density Functions (PDFs). The human annotations showed a sharp peak near class 0, indicating a highly consistent assignment of non-abusive class, despite disagreement, in the HS-Brexit dataset as seen in Figure 2. In contrast, the PDF for the strong perspectivist variant of the LLM showed a slightly right-skewed peak between 0.1 and 0.2, suggesting that the LLM assigned marginally higher soft labels than human annotators. The weak perspectivist PDF was flatter and more dispersed, with a small density spike near a probability of 0.2, reflecting greater uncertainty and inconsistency in labeling. The PDFs for the ConvAbuse dataset is presented in the Appendix D.

## 5.3 Prototypical *Persona* Annotators and Human Alignment

**Ablation 1**: Table 5 shows that annotator models trained on LLM annotations perform worse when tested on human labels, indicating a lack of alignment. The decline likely comes from the lack of corresponding match between LLM personas and the unknown individual human annotators.

**Ablation 2**: Figures 3 and 4 present an alignment analysis between LLM personas and human annotators. We compute cosine similarity between their annotation vectors. Using sample sizes of 5, 10, 50, and 100, stronger alignment was observed at sizes 50 and 100. In the ConvAbuse *strong* vari-

Figure 2: Figure showing the Probability Density Function illustrating Uncertainty in LLM annotations Vs Human in HSBrexit



| Human/LLM ID | Persona (strong) | Persona (weak) |
|---|---|---|
| 1 | Olivia, a female and white British person | white British female people |
| 2 | Emma, a female and white British person | white British with non-binary gender orientation |
| 3 | Ariel, a white person from the United Kingdom with a non-binary gender orientation | non-binary gender people from the United States |
| 4 | Sophia, a female and white person from the United Kingdom | white female people from the United Kingdom |
| 5 | Katrin, a female and white person from the United Kingdom | white female from United States |
| 6 | Eve, a female and white person from the United Kingdom | mixed Asian with a non-binary gender orientation |
| 7 | mixed Asian person | mixed Asian female |
| 8 | a white person from the United States with a non-binary gender orientation | mixed Asian female |

Figure 3: Figure showing Prototypical LLM annotators and Alignment with Human Annotators in ConvAbuse

ant, ANN(2–8) showed varying degrees of alignment with LLM Persona 1 (Olivia, female, white, British), while ANN(1) aligns more closely with LLM Persona 4 (Sophia, female, white, from the UK). Other LLM personas (2, 3, 5–8) exhibit no correspondence with any human annotator. We

further trained annotator models on annotations from LLM Personas 1 and 4, and evaluated them against human-labeled data. These models showed improved performance, approaching human-level results for both Composite Embedding and SBERT, as shown in Table 5.

| Model | SBERT | User Token | Composite Embedding | Composite Embedding + User Token |
|---|---|---|---|---|
| HL | 85.9 | 88.5 | 85.8 | 88.6 |
| LLM | 85.7 | 84.4 | 84.6 | 84.4 |
| LLM-H | 83.1 | 83.4 | 84.5 | 84.2 |
| LLM(1,4)-H | 85.4 | 82.6 | 85.1 | 85.9 |

Table 5: Model performance based on different training and testing label splits: HL (models trained and tested on Human Labels), LLM (models trained and tested on LLM Labels), LLM-H (models trained on LLM Labels, tested on Human Labels), and LLM(1,4)-H (models trained on the most aligned LLM personas 1 and 4 to human labels, tested on Human Labels).

In the HS-Brexit dataset, alignment is less consistent. In Figure 4, we see Persona 1, Male Muslim migrant, belonging to the *target* group mapped to annotators 4 and 5 of the human annotators belonging to the *control* group in the strong variant. Human annotators 1–3 belong to the Muslim or migrant group, while annotators 4–6 belong to the group with Western background, denoted as *locals*. Also, Persona 3 of the migrant group representing "neutral foreigner" shows positive alignment in the *weak* variant to the migrant group in human when "Muslim" was removed. These findings suggest that Llama2 includes prototypical personas capable of partially representing multiple human annotators. However, other defined personas fail to map to any observed human annotation patterns (cf. Appendix E).

## 6 Discussion and Conclusion

This work investigates Llama2's capacity to generate disaggregated labels for hate speech and offensiveness datasets using predefined personas, under two perspectivism frameworks: strong (individual) and weak (group) data perspectivism. We examine the quality and alignment of LLM-generated annotations with human-annotated datasets and evaluate downstream performance across existing annotator modeling techniques.

Llama2 annotations consistently exhibited higher inter-annotator agreement (Krippendorff's alpha ranging 0.55–0.91) than human annotations across both ConvAbuse and HS-Brexit datasets, though agreement decreased at higher temperatures. PDF analysis further indicated that LLM annotations tend to converge around features inherent in the model's underlying corpus, suggesting a divergence from human perspectives. As seen in Figure 2, the PDF using the soft label distribution of the abusive class shows human annotations aligning towards the non-abusive class, strong perspectivism aligning more towards the abusive class, and

weak perspectivism showing a relatively flat and dispersed distribution depicting high uncertainty.

In terms of performance of annotator modeling methods, LLM annotations shifted model efficacy. While prior work confirmed that annotator models trained on human-annotated datasets with high agreement (e.g., ConvAbuse) performed best with the Composite Embedding + User Token model, and those with low agreement (e.g., HS-Brexit) favored the User Token model, our findings with LLM-generated annotations demonstrate that simpler models, specifically SBERT and Composite Embedding models without explicit annotator information, showed improved results. This shift implies that LLM-generated annotations align more with generalized perspectives and are less suited to highly personalized approaches. Comparing the two perspectivism approaches, strong data perspectivism on ConvAbuse, characterized by overlapping and more personalized features, improved the performance of annotator modeling techniques over its weak counterpart. Conversely, weak perspectivism on HS-Brexit, with its contrasting demographic features in groups, yielded improved performance specifically with SBERT and Composite Embedding models, suggesting that contrasting demographic diversity tends to influence the choice of perspectivism approach and annotator modeling performance in LLMs.

Our ablation studies revealed LLM personas do not directly correspond to human annotators. However, as seen in Figure 3, we identified generalized "prototypical persona features" working as representatives of groups of humans (e.g., ANN 2-8 mapping to LLM Persona 1, ANN1 to LLM Persona 4). Swapping the labels of corresponding annotators in the original dataset with these prototypical annotator labels, and evaluating with the human test set, slightly improved results, as seen in Table 5, presenting a novel approach for modeling perspectivism in LLMs. These findings suggest that while

127

Figure 4: Figure showing Prototypical LLM annotators and Alignment with Human Annotators in HSBrexit

LLMs offer insights into subjective domains, their capacity to fully embody external personas remains limited to their underlying corpus, supporting an aggregated view rather than personalization. Future work should focus on standardization and generate more diversified personas, systematically varying features, and expanding evaluation to other LLMs to fully investigate these prototypical attributes and their potential in capturing a wider scope of perspectives.

## 7 Limitations

This study is based on two datasets and focuses exclusively on binary classification tasks for hate and offensive speech detection. One potential limitation is that the data used to train Llama2-13B may have been filtered, reducing its sensitivity to detecting abusive content, potentially influencing the observed results. Our analysis is also limited to this model, and we did not investigate how newer variants of Llama or other LLMs, like GPT 4o, might influence the results. The personas used for generating annotations were limited to the demographic features explicitly provided in the original datasets, with slight modifications to fit the perspectivist spectrum. Furthermore, we did not quantify the extent to which the model's attention was distributed between the persona and the input sentences. Understanding this balance could provide deeper insight into how strongly LLMs personalize their annotations.

Another limitation of this study arises from the design of the annotation prompt for the HS-Brexit dataset variant, which focused on 'hate speech towards Brexit'. However, the prompt was structured to provide general contextual information about Brexit and simulate the prior knowledge of human annotators. A follow-up experiment analysing the

model's attention mechanism revealed that the instance of "Brexit" appearing first in the prompt received a significantly higher attention score of 0.0654 than the "Brexit" label target which received an attention score of 0.0050. Furthermore, when 'immigrants' was targeted instead, it received an attention score of 0.0117, which was higher than that given to 'Brexit' as a target. This suggests that the models have learned to recognise plausible targets for hate speech, which warrants further investigation. However, this paper's specific focus is to investigate the impact of Annotator Personas on LLM behaviour across the perspectivism spectrum. It therefore does not include a deep analysis of the model's sensitivity to target plausibility. Nevertheless, we present this as a compelling avenue for future research, while maintaining that our core findings regarding persona-driven perspectivism remain valid within the described experimental setup. Our codes are publicly available[1] to support future work.

## Acknowledgments

## References

Gati V. Aher, Rosa I. Arriaga, and Adam Tauman Kalai. 2023. Using large language models to simulate multiple humans and replicate human subject studies. In *International Conference on Machine Learning, ICML 2023*, pages 337–371. PMLR.

Sohail Akhtar, Valerio Basile, and Viviana Patti. 2021. Whose opinions matter? perspective-aware models to identify opinions of hate speech victims in abusive language detection. *Preprint*, arXiv:2106.15896.

Dina Almanea and Massimo Poesio. 2022. Armis - the arabic misogyny and sexism corpus with annotator subjective disagreements. In *Proceedings of the Thirteenth Language Resources and Evaluation Conference*, pages 2282–2291. European Language Resources Association.

Lora Aroyo and Chris Welty. 2015. Truth is a lie: Crowd truth and the seven myths of human annotation. *AI Magazine*, 36(1):15–24.

Anna Bavaresco, Raffaella Bernardi, Leonardo Bertolazzi, Desmond Elliott, Raquel Fernández, Albert

Gatt, Esam Ghaleb, Mario Giulianelli, Michael Hanna, Alexander Koller, André F. T. Martins, Philipp Mondorf, Vera Neplenbroek, Sandro Pezzelle, Barbara Plank, David Schlangen, Alessandro Suglia, Aditya K Surikuchi, Ece Takmaz, and Alberto Testoni. 2024. Llms instead of human judges? a large scale empirical study across 20 nlp evaluation tasks. *Preprint*, arXiv:2406.18403.

Tilman Beck, Hendrik Schuff, Anne Lauscher, and Iryna Gurevych. 2024. Sensitivity, performance, robustness: Deconstructing the effect of sociodemographic prompting. In *Proceedings of the 18th Conference of the European Chapter of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 2589–2615. Association for Computational Linguistics.

Federico Cabitza, Andrea Campagner, and Valerio Basile. 2023. Toward a perspectivist turn in ground truthing for predictive computing. In *Thirty-Seventh AAAI Conference on Artificial Intelligence*, pages 6860–6868. AAAI Press.

Amanda Cercas Curry, Gavin Abercrombie, and Verena Rieser. 2021. Convabuse: Data, analysis, and benchmarks for nuanced abuse detection in conversational ai. In *Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing*, pages 7388–7403. Association for Computational Linguistics.

Xia Cui. 2023. xiacui at semeval-2023 task 11: Learning a model in mixed-annotator datasets using annotator ranking scores as training weights. In *Proceedings of the 17th International Workshop on Semantic Evaluation (SemEval-2023)*, pages 1076–1084. Association for Computational Linguistics.

Aida Mostafazadeh Davani, Mark Díaz, and Vinodkumar Prabhakaran. 2022. Dealing with disagreements: Looking beyond the majority vote in subjective annotations. *Transactions of the Association for Computational Linguistics*, 10:92–110.

Dorottya Demszky, Dana Movshovitz-Attias, Jeongwoo Ko, Alan Cowen, Gaurav Nemade, and Sujith Ravi. 2020. Goemotions: A dataset of fine-grained emotions. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 4040–4054. Association for Computational Linguistics.

Eve Fleisig, Su Lin Blodgett, Dan Klein, and Zeerak Talat. 2024. The perspectivist paradigm shift: Assumptions and challenges of capturing human labels. In *Proceedings of the 2024 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies (Volume 1: Long Papers)*, pages 2279–2292. Association for Computational Linguistics.

Simona Frenda, Gavin Abercrombie, Valerio Basile, Alessandro Pedrani, Raffaella Panizzon, Alessandra Teresa Cignarella, Cristina Marco, and Davide

[1] https://doi.org/10.5281/zenodo.16744588

Bernardi. 2024. Perspectivist approaches to natural language processing: a survey. *Language Resources and Evaluation*.

José Antonio García-Díaz, Ronghao Pan, Gema Alcaráz-Mármol, María José Marín-Pérez, and Rafael Valencia-García. 2023. Umuteam at semeval-2023 task 11: Ensemble learning applied to binary supervised classifiers with disagreements. In *Proceedings of the 17th International Workshop on Semantic Evaluation (SemEval-2023)*, pages 1061–1066. Association for Computational Linguistics.

Fabrizio Gilardi, Meysam Alizadeh, and Maël Kubli. 2023. Chatgpt outperforms crowd workers for text-annotation tasks. *Proceedings of the National Academy of Sciences*, 120(30):e2305016120.

Kristina Gligorić, Tijana Zrnic, Cinoo Lee, Emmanuel J. Candès, and Dan Jurafsky. 2024. Can unconfident llm annotations be used for confident conclusions? *Preprint*, arXiv:2408.15204.

Shirley Anugrah Hayati, Minhwa Lee, Dheeraj Rajagopal, and Dongyeop Kang. 2024. How far can we extract diverse perspectives from large language models? In *Proceedings of the 2024 Conference on Empirical Methods in Natural Language Processing*, pages 5336–5366. Association for Computational Linguistics.

Tiancheng Hu and Nigel Collier. 2024. Quantifying the persona effect in llm simulations. *Preprint*, arXiv:2402.10811.

Olivia Huang, Eve Fleisig, and Dan Klein. 2023. Incorporating worker perspectives into mturk annotation practices for nlp. In *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing*, pages 1010–1028. Association for Computational Linguistics.

Jonathan Ivey, Shivani Kumar, Jiayu Liu, Hua Shen, Sushrita Rakshit, Rohan Raju, Haotian Zhang, Aparna Ananthasubramaniam, Junghwan Kim, Bowen Yi, Dustin Wright, Abraham Israeli, Anders Giovanni Møller, Lechen Zhang, and David Jurgens. 2024. Real or robotic? assessing whether llms accurately simulate qualities of human responses in dialogue. *Preprint*, arXiv:2409.08330.

Brendan Kennedy, Mohammad Atari, Aida Mostafazadeh Davani, Leigh Yeh, Ali Omrani, Yehsong Kim, Kris Coombs, Shreya Havaldar, Gwenyth Portillo-Wightman, Elaine Gonzalez, Joe Hoover, Aida Azatian, Alyzeh Hussain, Austin Lara, Gabriel Cardenas, Adam Omary, Christina Park, Xin Wang, Clarisa Wijaya, Yong Zhang, Beth Meyerowitz, and Morteza Dehghani. 2022. Introducing the gab hate corpus: defining and applying hate-based rhetoric to social media posts at scale. *Language Resources and Evaluation*, 56(1):79–108.

Matej Klemen and Marko Robnik-Šikonja. 2022. Ulfri at semeval-2022 task 4: Leveraging uncertainty and additional knowledge for patronizing and condescending language detection. In *Proceedings of the 16th International Workshop on Semantic Evaluation (SemEval-2022)*, pages 525–532. Association for Computational Linguistics.

Elisa Leonardelli, Gavin Abercrombie, Dina Almanea, Valerio Basile, Tommaso Fornaciari, Barbara Plank, Verena Rieser, Alexandra Uma, and Massimo Poesio. 2023. Semeval-2023 task 11: Learning with disagreements (lewidi). In *Proceedings of the 17th International Workshop on Semantic Evaluation (SemEval-2023)*, pages 2304–2318. Association for Computational Linguistics.

Sarumi Oluyemi, Béla Neuendorf, Joan Plepi, Lucie Flek, Jörg Schlötterer, and Charles Welch. 2024. Corpus considerations for annotator modeling and scaling. In *Proceedings of the 2024 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies (Volume 1: Long Papers)*, pages 1029–1040. Association for Computational Linguistics.

Maja Pavlovic and Massimo Poesio. 2024a. The effectiveness of llms as annotators: A comparative overview and empirical analysis of direct representation. In *Proceedings of the 3rd Workshop on Perspectivist Approaches to NLP (NLPerspectives) @ LREC-COLING 2024*, pages 100–110. ELRA and ICCL.

Maja Pavlovic and Massimo Poesio. 2024b. Understanding the effect of temperature on alignment with human opinions. *Preprint*, arXiv:2411.10080.

Joan Plepi, Béla Neuendorf, Lucie Flek, and Charles Welch. 2022. Unifying data perspectivism and personalization: An application to social norms. In *Proceedings of the 2022 Conference on Empirical Methods in Natural Language Processing*, pages 7391–7402. Association for Computational Linguistics.

Paul Rottger, Bertie Vidgen, Dirk Hovy, and Janet Pierrehumbert. 2022. Two contrasting data annotation paradigms for subjective nlp tasks. In *Proceedings of the 2022 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pages 175–190. Association for Computational Linguistics.

Michael Sullivan, Mohammed Yasin, and Cassandra L. Jacobs. 2023. University at buffalo at semeval-2023 task 11: Masda–modelling annotator sensibilities through disaggregation. In *Proceedings of the 17th International Workshop on Semantic Evaluation (SemEval-2023)*, pages 978–985. Association for Computational Linguistics.

Huaman Sun, Jiaxin Pei, Minje Choi, and David Jurgens. 2025. Sociodemographic prompting is not yet an effective approach for simulating subjective judgments with llms. *Preprint*, arXiv:2311.09730.

Nikolas Vitsakis, Amit Parekh, Tanvi Dinkar, Gavin Abercrombie, Ioannis Konstas, and Verena Rieser.

2023. ilab at semeval-2023 task 11 le-wi-di: Modelling disagreement or modelling perspectives? In *Proceedings of the 17th International Workshop on Semantic Evaluation (SemEval-2023)*, pages 1660–1669. Association for Computational Linguistics.

Angelina Wang, Jamie Morgenstern, and John P. Dickerson. 2024. Large language models should not replace human participants because they can misportray and flatten identity groups. *Preprint*, arXiv:2402.01908.

Jin Xu, Mariët Theune, and Daniel Braun. 2024. Leveraging annotator disagreement for text classification. In *Proceedings of the 7th International Conference on Natural Language and Speech Processing (IC-NLSP 2024)*, pages 1–10. Association for Computational Linguistics.

# Appendix

## A   Prompt Design

**Prompt for Strong HS-Brexit** *annotate this statement regarding Brexit. As a (persona) in the UK, you must assign a label to the statement. Use 'O' if the statement does not express hate towards Brexit, and '1' if it does express hate towards Brexit. Provide only the number as your response without any additional text or explanation for example "###Annotator:0" or "###Annotator:1"*

**Prompt for Weak HS-Brexit** *annotate this statement regarding Brexit.  for a group of (persona) in the UK or Europe, you must assign a label to the statement. Use 'O' if the statement does not express hate towards Brexit, and '1' if it does express hate towards Brexit. Provide only the number as your response without any additional text or explanation for example "### Annotator:0" or "###Annotator:1"*

**Prompt for Weak ConvAbuse** *annotate these conversations between users and bots. As a group of persona within the age range of 19-21, a level 1 English speaker and an Undergraduate student in Gender Studies and Sociology, you must assign a label to the statement.  Use 'O' if the statement is not abusive, and '1' if it does express abusiveness by the user towards the bot or agent. Provide only the number as your response without any additional text or explanation.  for example "### Annotator:0" or "### Annotator:]"*

**Prompt for Strong ConvAbuse** *annotate these conversations between users and bots. As a persona within the age range of 19-21, a level 1 English speaker and an Undergraduate student in Gender Studies and Sociology, you must assign a label to the statement. Use 'O' if the statement is not abusive, and '1' if it does express abusiveness by the user towards the bot or agent.  Provide only the number as your response without any additional text or explanation. for example "### Annotator:0" or "### Annotator:]"*

## B   Persona Descriptions

### HS-Brexit Persona for Strong Perspectives

- Male Muslim Migrant

- Female Muslim Migrant

- Neutral foreigner in the UK

- Native English man

- Native English Woman

- Liberal English person

### HS-Brexit Persona for Weak Perspectives

- researchers with Western background having experience in linguistic annotation

- first or second generation muslim immigrant students from developing countries

### ConvAbuse Persona for Weak Perspectives

- white British female people

- white British with non-binary gender orientation

- non-binary gender people from the United States

- white female people from the United Kingdom

- white female from United States

- mixed Asian with a non-binary gender orientation

- mixed Asian female

- white people from the United States with a non-binary gender orientation

### ConvAbuse Persona for Strong Perspectives

- Olivia, a female and white british person

- Emma, a female and white british person

- Ariel, a white person from the United Kingdom with a non-binary gender orientation

- Sophia, a female and white person from the United Kingdom

- Katrin, a female and white person from the United Kingdom

- Eve, a female and white person from the United Kingdom

- a mixed Asian person

- a white person from the United States with a non-binary gender orientation

# C   Model performance for Strong and Weak Data Perspectivism

| Model | $\alpha$ | User-Token | Composite | Composite+ User-Token | Multitasking | SBERT |
|---|---|---|---|---|---|---|
| **Strong Perspectivism** | | | | | | |
| Human | 0.65 | 88.5 | 85.8 | 88.6 | 82.3 | 85.9 |
| 0 | 0.91 | 84.1 | 83.0 | 84.4 | 46.9 | 83.1 |
| 0.1 | 0.87 | 84.4 | **84.6** | 84.3 | 81.1 | **85.7** |
| 0.2 | 0.81 | 80.5 | 81.5 | 80.0 | 46.8 | 81.5 |
| 0.5 | 0.68 | 69.9 | 70.7 | 71.1 | 45.1 | 69.4 |
| 0.8 | 0.60 | 63.5 | 65.1 | 64.4 | 62.6 | 64.6 |
| **Weak Perspectivism** | | | | | | |
| 0 | 0.93 | 83.7 | 83.7 | 81.8 | 69.0 | 85.2 |
| 0.1 | 0.88 | 80.1 | 79.3 | 78.3 | 79.8 | 82.0 |
| 0.2 | 0.82 | 81.2 | 81.5 | 81.2 | 70.3 | 82.1 |
| 0.5 | 0.67 | 71.7 | 69.7 | 71.4 | 64.7 | 69.5 |
| 0.8 | 0.62 | 61.7 | 61.4 | 62.3 | 58.1 | 61.4 |

Table 6: Performance of Annotator modeling methods for Strong and Weak data Perspectivism (**ConvAbuse dataset**) across various temperatures.

| Model | $\alpha$ | User-Token | Composite | Composite+ User-Token | Multitasking | SBERT |
|---|---|---|---|---|---|---|
| **Strong Perspectivism** | | | | | | |
| Human | 0.35 | 77.6 | 67.6 | 77.3 | 71.7 | 68.6 |
| 0 | 0.81 | 69.3 | 71.3 | 71.2 | 65.1 | 72.2 |
| 0.1 | 0.73 | 69.4 | 71.8 | 71.0 | 61.8 | 69.2 |
| 0.2 | 0.67 | 66.3 | 63.8 | 61.9 | 61.4 | 67.2 |
| 0.5 | 0.62 | 61.5 | 61.3 | 61.4 | 49.5 | 62.2 |
| 0.8 | 0.58 | 52.4 | 56.1 | 54.2 | 51.2 | 56.6 |
| **Weak Perspectivism** | | | | | | |
| 0 | 0.75 | 72.2 | **72.4** | 71.7 | 60.3 | **73.2** |
| 0.1 | 0.69 | 66.6 | 65.8 | 65.5 | 62.0 | 69.1 |
| 0.2 | 0.62 | 62.2 | 63.8 | 69.9 | 59.2 | 66.8 |
| 0.5 | 0.54 | 58.0 | 58.4 | 57.9 | 39.2 | 56.1 |
| 0.8 | 0.55 | 55.2 | 57.8 | 56.7 | 55.4 | 56.6 |

Table 7: Performance of Annotator modeling methods for Strong and Weak data Perspectivism (**HS-Brexit dataset**) across various temperatures.

# D   Probability Density Function for Uncertainty and Annotation Quality

The Figure 5 shows the probability density function of the weak data perspectivism in ConvAbuse using the majority class as a reference point.



Figure 5: Probability Density Function ConvAbuse Dataset

# E   Human Vs Persona Alignment and Prototypes ConvAbuse Dataset

| Human Annotator | SampleSize | Best Match LLM Persona | Similarity score |
|---|---|---|---|
| 2 | 100 | 1 | 0.791 |
| 3 | 100 | 1 | 0.894 |
| 7 | 100 | 1 | 0.913 |
| 8 | 100 | 1 | 0.671 |
| 1 | 100 | 4 | 0.707 |
| 6 | 100 | 1 | 0.707 |
| 4 | 100 | 1 | 0.816 |
| 5 | 100 | 1 | 0.745 |

Table 8: Mapping of Human Annotators to Best Matching LLM Personas based on Cosine Similarity.

**Prototypical Annotators and their Alignment with Human Annotators Across Varying Sample sizes**

Alignment between Human Annotator and Best Matching LLM in strong data perspectivism-HSBrexit



Figure 6: Showing the identified Prototypical annotators in HS-Brexit dataset and the alignment with human annotators

Figure 7: Showing the identified Prototypical annotators in ConvAbuse dataset and the alignment with human annotators

# The Chunking Paradigm: Recursive Semantic for RAG Optimization

**Seemab Latif** [*]**, Huma Ameer**[†]**, Muhammad Hannan Akram**[†] **and Mehwish Fatima**[†]
School of Electrical Engineering and Computer Science,
National University of Sciences and Technology (NUST), Islamabad, Pakistan
{seemab.latif, hameer.msds20seecs,makram.bsds23seecs,
mehwish.fatima}@seecs.edu.pk

## Abstract

Retrieval Augmented Generation (RAG) has risen to prominence for boosting the capabilities of Large Language Models (LLMs) through the integration of external knowledge. Notably, the document chunking process plays a central role in the performance of RAG pipelines. Nevertheless, incoherent document splits and inappropriate chunk sizes hinder retrieval efficiency and contextual accuracy. To address this, we propose Recursive Semantic Chunking (RSC), a dynamic and adaptive chunking framework that ensures semantic coherence. It maintains coherence by recursively splitting large chunks and merging smaller ones. Unlike conventional methods, RSC preserves contextual integrity while optimizing retrieval efficiency. The evaluation across 4 distinct datasets outperformed traditional semantic chunking techniques on evaluation metrics; contextual relevancy, contextual precision, contextual recall, retrieval time, faithfulness and answer relevancy. Results demonstrate that RSC consistently outperforms traditional chunking techniques, achieving higher contextual relevancy and total score while maintaining efficient retrieval times. These findings highlight the potential to optimize RAG systems and to improve the document chunking steps in the systems.

## 1 Introduction

Large Language Models (LLMs) are widely adopted across various domains in the form of chatbots, AI assistants, and other applications (Siddharth and Luo, 2024; Sahlman et al., 2023). The performance of LLMs is enhanced via the integration of external knowledge sources, specifically for custom applications. In addition, we can leverage the capabilities of LLMs without training them. The aforementioned enhancement can be made via

Retreival-Augmented Generation (RAG) (Lewis et al., 2020).

The RAG process begins with a user's query being sent to the LLM, which generates a retrieval request based on that query. This request is forwarded to the retriever system, which searches the vector database. Embeddings of documents chunk i.e. context is stored in vector database. The relevant context is then retrieved and combined with the user's query before being sent to the LLM for a final response, as shown in Figure 1. Researchers have developed various RAG-based solutions across different domains, such as finance and healthcare (Alkhalaf et al., 2024; He et al., 2024; Feng et al., 2024; Mathur et al., 2024).

The critical aspect of the RAG pipeline is the chunking of documents. Chunking in RAG systems is a technique that breaks down large documents into smaller, manageable segments known as "chunks" (LangChain, 2024). This process is crucial as it enhances the efficiency and accuracy of information retrieval, which leads to better outcomes for the system. The nature of context retrieved from the vector database is based on the segmentation of these documents, therefore, the choice of chunking techniques is a significant step in the pipeline (Setty et al., 2024). The chunking techniques directly affect the quality of retrieved-context and retrieval time. It eventually affects the quality of the product that is utilizing RAG-based applications. The choice of chunking is quite challenging i.e. larger chunks can lead to slower retrieval, or retrieve irrelevant chunks and small chunks may not adhere to a coherent information unit. Recently, there has been a shift in research focus towards optimal chunking techniques i.e. (Yepes et al., 2024). Although frameworks such as LangChain (AI, 2024) and LamaIndex (Liu, 2022) have various chunking strategies. Due to complexities of the document structure, and cus-

---

[*]Corresponding author.
[†]Equal contribution.

Figure 1: Information Flow in Retrieval Augmented Generation (RAG)

tom systems, it is still a challenging task at hand.

In this paper, we propose Recursive Semantic Chunking that focuses on optimizing the semantic chunking of the documents. The following are the key contributions of this work:

- We propose a Recursive Semantic Chunking method, designed to split textual documents into coherent semantic chunks of an appropriate size.

- We introduce a dynamic method for adjusting the chunk size. The recursive nature of this method ensures that larger text segments are progressively broken down while maintaining semantic integrity. Smaller segments are merged strategically, keeping the chunk size balanced.

- We demonstrate that the proposed method improves retrieval time compared to traditional chunking techniques.

- As part of this work, we introduce NewsMatrix-71, a large-scale, multi-domain news dataset.

## 2   Related Work

Retrieval Augmented Generation systems rely on the context returned from the retrieval algorithms, making chunking a key factor in the pipeline (Yepes et al., 2024). Therefore, the choice of chunking strategies is a critical step. Ineffective techniques can result in either incomplete chunks leading to losing context or large chunks with irrelevant information negatively impacting the accuracy of the retrieval (Setty et al., 2024).

One of the common approaches is to split the document based on fixed numbers of chunks. However, it has a potential loss of context in both cases larger or smaller chunk size (Teja, 2023). To address this, the researchers introduce recursive split by character technique (LangChain, 2023). It recursively splits keeping the longest text chunks together with a need to define and constant adjustment of chunk size overlapping making it computationally expensive.

Although the recursive text split tends to keep the chunks semantically closed together, it does not directly account for semantic meaning. Conversely, semantic chunking (LangChain, 2024) groups the text that is semantically similar together. It first splits the text into sentences and groups them into three sentences, then merges similar groups in the embedding space. However, this technique does not ensure optimal chunk sizes. Since its mechanism is dependent on the similarity of the embedding vectors, it may lead to larger chunks and cause hallucinations.

Agentic chunking (FullStackRetrieval, 2024) pushed this idea further by leveraging Large Language Models. It converts text into propositions via LLMs (Chen et al., 2024). Propositions are defined as standalone statements that convey a single fact clearly without needing extra context. It can be referred to as the smallest unit of meaning within a text, each expressing one distinct idea. Propositions retain the semantic meaning in individual statements as shown in the following example:

Once the propositions are created, these are passed to an LLM, which is then prompted to group these chunks. This approach offers flexibility and higher accuracy. Nevertheless, it requires well-crafted prompts and dependency on the capability of the acquired LLM.

Working on efficient chunking techniques is an open research area as not much has been explored in this regard.

138

**Proposition Conversion Example**

**Original Sentence**

*"Three new products were launched this year, expanding our reach into international markets."*

**Converted Propositions**

*"Three new products were launched this year."*
*"The company expanded into international markets."*

## 3 Dataset

We introduce a new large-scale news dataset, named NewsMatrix-71[1]. It covers a diverse set of news categories over multiple years.

### 3.1 Scraped News Dataset (NewsMatrix-71)

We compile this dataset by scraping English news articles from Dawn[2], Tribune[3], and Daily Times[4]. This dataset covers the span of three years (2021–2023) and has up to 96,859 news articles categorized into 71 unique topics, including Business, Fashion, Health, World, and more. It offers a diverse, time-spanning, and category-rich corpus suitable for various NLP tasks. It captures a broad spectrum of global and regional news, making it a valuable resource for research. Given the size and scope of this dataset, we will selectively release a publicly available subset to facilitate reproducibility and further research.

## 4 Recursive Semantic Chunking

This section presents the Recursive Semantic Chunking framework in detail. The primary objective is to ensure the splitting of chunks is semantically coherent and maintains the integrity of the content. In addition, the size of the chunks should be optimal. The standard semantic chunking technique tends to generate large chunks, which

---

[1]This data will be published publicly and free for research purposes after the paper's acceptance. It will be shared under the **Creative Commons Attribution 4.0 International License (CC BY 4.0)**

[2]Dawn

[3]Tribune

[4]Daily Times

negatively impact the performance of retrieval-augmented generation systems. Furthermore, in custom RAG projects, documents often belong to specific topics, and larger chunks reduce system efficiency.

Algorithm 1 provides a detailed outline of the proposed chunking process. All predefined values are determined after extensive experimentation. The following steps describe the pipeline.

**Segmentation of Textual Data from Files**

The data store consists of files $f_i$ containing textual data stored as strings $T_i$. Since LLMs have token limits, each $T_i$ undergoes a length check. If it exceeds the threshold $T_{\max}$, the file is split into smaller segments $\{t_1, t_2, \ldots, t_n\}$, ensuring that $|t_j| \leq T_{\max}$. The splitting occurs at the nearest sentence boundary (e.g., full stop, question mark) to preserve linguistic coherence.

**Initial Semantic Chunking**

Each segment $t_j$ undergoes an initial semantic chunking process (LangChain, 2024). In this step, the semantically similar texts are grouped in the embedding space, forming $C_0 = \{c_1, c_2, \ldots, c_m\}$, where $c_k$ represents an initial chunk.

**Recursive Semantic Chunking**

For each chunk $c_k \in C_0$, the semantic chunker is recursively applied if its length exceeds the threshold $T_{\text{chunk}}$ (1,500 characters). With each recursive iteration, the breakpoint threshold parameter is gradually reduced, ensuring that large chunks are broken into smaller, semantically meaningful segments. The recursive function $R(c, T)$ operates as follows:

$$R(c, T) = \begin{cases} c & \text{if } |c| \leq T \\ R(\text{split}(c, T - \delta), T - \delta) & \text{if } |c| > T \end{cases}$$

where $\delta$ represents a small reduction factor to progressively decrease chunk size in each iteration. The reduction factor $\delta$ is heuristically set to 3 after initial experimentation. Although not tuned through systematic search, this value is selected to ensure a gradual and controlled recursive breakdown of large chunks. This value is kept fixed across all datasets to maintain consistency and reproducibility.

## Merging Short Chunks

Following recursive chunking, some chunks may become too short (i.e., less than $T_{\text{merge}}$, set to 350 characters). Extremely small chunks may lack semantic coherence, leading to information loss. To address this, the similarity score of smaller chunks is calculated with both preceding and subsequent chunks. It is merged with the chunk that has the highest similarity score. This ensures semantic integrity while preventing the loss of meaningful text. The merging process is defined as follows:

$$
\text{For } i = 1 \text{ to } n : \begin{cases} \text{If } |c_i| < T_{\text{merge}}, \text{ compute:} \\ \quad S_{\text{prev}} = \text{sim}(c_i, c_{i-1}) \\ \quad S_{\text{next}} = \text{sim}(c_i, c_{i+1}) \\ \text{Merge with highest similarity chunk} \\ \quad \text{If } S_{\text{prev}} \geq S_{\text{next}}, \text{ then:} \\ \qquad c_{i-1} \leftarrow c_{i-1} + c_i \\ \quad \text{Else:} \\ \qquad c_{i+1} \leftarrow c_i + c_{i+1} \end{cases}
$$

Here, $S_{\text{prev}}$ and $S_{\text{next}}$ represent the similarity scores between the small chunk $c_i$ and its neighboring chunks $c_{i-1}$ and $c_{i+1}$, respectively. The chunk $c_i$ is merged with the chunk that has the highest similarity score, ensuring that the resulting merged chunk maintains semantic coherence.

## Uniform Chunk Size Adjustment

Finally, the algorithm checks whether any chunk exceeds the threshold $T_{\text{final}}$ (2,500 characters). If a chunk surpasses this limit, it undergoes a recursive character-based text split (LangChain, 2023). The final adjustment process is defined as:

$$
\text{For } i = 1 \text{ to } m : \begin{cases} \text{If } |c_i| > T_{\text{final}} : \\ \quad \text{Apply Recursive Split Function:} \\ \quad c_i \leftarrow \text{RecursiveSplit}(c_i, T_{\text{final}}) \end{cases}
$$

This step ensures that the final chunk set, $C_{\text{final}} = \{c_1, c_2, \ldots, c_m\}$, meets size constraints while maintaining semantic coherence. The processed chunks are then stored in vector databases for RAG tasks.

## Distinction from Baseline Chunkers

While our method incorporates components from existing LangChain utilities, i.e. semantic chunking for initial grouping and character-based recursive splitting for final chunk size enforcement. These steps function as structural helpers rather than the core of our approach. The key innovation of RSC lies in its intermediate refinement

---

**Algorithm 1:** Recursive Semantic Chunking

**Input:** Dataset $D = \{f_1, f_2, \ldots, f_N\}$;
Maximum chunk size $T_{\text{max}} = 15{,}000$;
Recursive chunking threshold $T_{\text{chunk}} = 1{,}500$;
Final chunk size threshold $T_{\text{final}} = 2{,}500$;
Minimum chunk size for merging $T_{\text{merge}} = 350$
**Output:** Final set of chunks $C_{\text{final}}$

1  **Initialization:**
2  $C_{\text{final}} \leftarrow \emptyset$
3  **Initial Semantic Chunking:**
4  Apply chunking to each segment $t_j$:
5  $C_0 \leftarrow \{c_1, c_2, \ldots, c_m\}$
6  **foreach** *chunk* $c_k \in C_0$ **do**
7  $\quad$ **if** $|c_k| > T_{chunk}$ **then**
8  $\qquad$ **Recursive Semantic Chunking:**
9  $\qquad R(c_k, T_{\text{chunk}}) =$
$\qquad\quad R(\text{split}(c_k, T_{\text{chunk}} - \delta), T_{\text{chunk}} - \delta)$
10 $\qquad c_k \leftarrow R(c_k, T_{\text{chunk}})$

11 **foreach** *chunk* $c_k \in C_0$ **do**
12 $\quad$ **if** $|c_k| \leq T_{merge}$ **then**
13 $\qquad$ Compute similarity with previous chunk:
14 $\qquad S_{\text{prev}} \leftarrow \text{similarity}(c_{k-1}, c_k)$
15 $\qquad$ Compute similarity with next chunk:
16 $\qquad S_{\text{next}} \leftarrow \text{similarity}(c_k, c_{k+1})$
17 $\qquad$ **if** $S_{prev} \geq S_{next}$ **then**
18 $\qquad\quad$ Merge with previous chunk:
19 $\qquad\quad c_{k-1} \leftarrow c_{k-1} + c_k$
20 $\qquad$ **else**
21 $\qquad\quad$ Merge with next chunk:
22 $\qquad\quad c_{k+1} \leftarrow c_k + c_{k+1}$

23 Add merged chunks to $C_{\text{final}}$
24 **foreach** *chunk* $c_k \in C_{final}$ **do**
25 $\quad$ **if** $|c_k| > T_{final}$ **then**
26 $\qquad$ Split chunk:
27 $\qquad c_k \leftarrow \text{split}(c_k, T_{\text{final}})$

28 **Return:** Final set of chunks $C_{\text{final}}$

---

logic: recursive breakdown with dynamic thresholds, similarity-based merging of smaller chunks, and controlled preservation of semantic coherence. These operations are not present in the baseline LangChain chunkers and are designed to address the limitations of fixed-size or purely embedding-based segmentation. Therefore, while we leverage LangChain for low-level chunk initialization and splitting, the significant performance improvements observed in contextual and answer-level metrics stem from our recursive and adaptive chunking strategy.

## 5   Experimental Design

Our evaluation framework is designed to rigorously assess the impact of our proposed technique: Recursive Semantic Chunking (RSC). Incorporating RSC in the RAG pipeline for question-answering tasks, we demonstrate its capabilities in preserving contextual coherence and improving retrieval precision. This section details our evaluation methodology, covering dataset selection, synthetic question

Table 1: Summary of Datasets used for Evaluating the Proposed Chunking Technique, including Open-source Corpora and the Custom Dataset NewsMatrix-71.

| Dataset | Words | Characters | Paragraphs | Source |
|---|---|---|---|---|
| BBC | 854,490 | 5,039,982 | 2,225 | BBC Dataset |
| SQuAD | 152,394 | 966,345 | 1,204 | SQuAD |
| QuaC | 440,971 | 2,664,801 | 1,000 | QuaC |
| NewsMatrix-71 | 677,258 | 4,227,679 | 1,500 | Dawn, Tribune Daily Times |

generation, chunking techniques, implementation setup, and performance metrics

## 5.1 Datasets

We evaluate our proposed chunking technique using four datasets, including three open-source corpora—BBC (Greene and Cunningham, 2006), SQuAD (Rajpurkar et al., 2016), and QuaC (Choi et al., 2018)—along with a custom-scraped news dataset, NewsMatrix-71. The NewsMatrix-71 dataset, created by scraping English news articles, is stored in .txt format. For experimentation, we use a 1,500-article subset containing 677,258 words and 4,227,679 characters. A summary of all datasets is provided in Table 1.

## 5.2 Synthetic Question Generation

These evaluations of the chunking techniques are based on the response from the question-answering system. Therefore, we utilized LLM to create synthetic questions from each dataset. For each dataset, we randomly generate 50 synthetic questions per dataset to balance computational feasibility with evaluation diversity. This quantity is consistent with recent study Merola and Singh, 2025. This quantity is consistent with recent study Merola and Singh, 2025. To generate synthetic questions, we randomly selected passages from each dataset. To ensure reasonable topic coverage, we manually examined multiple random subsets and selected one for question generation. While this approach does not guarantee perfect topic stratification, it provides a practical balance between topic diversity and simplicity in sampling. We employ Gemini Flash 1.5 to generate corresponding questions. The ChatPromptTemplate module from LangChain is used to structure the input prompt, guiding the model to generate relevant and context-aware questions for each passage. Once generated, the synthetic questions are stored and later used to assess the retrieval and response quality of different chunking techniques. By introducing synthetic

queries, we create an additional layer of evaluation that allows us to measure how well-chunked text segments support question-answering tasks beyond the scope of existing datasets.

## 5.3 Chunking Techniques

To establish a baseline, we implement three widely used chunking techniques. Recursive Character Text Splitter segments (LangChain, 2023), and Semantic Chunking (LangChain, 2024). Next, we employ our proposed technique; Recursive Semantic Chunking framework for comparison.

## 5.4 Implementation Details

For downstream question-answering tasks, we store the chunks in the RAG pipeline using LangChain[5]. All the techniques use "*all-MiniLM-L6-v2*" [6]embedding. The resulting chunks are stored in the Facebook AI Similarity Search (FAISS) vector database (Douze et al., 2024). The "*ChatPromptTemplate module*" is used with "*Gemini Flash 1.5*" [7], a state-of-the-art Large Language Model optimized for contextual reasoning.

## 5.5 Evaluation metrics

We assess chunking techniques by integrating them into the RAG pipeline for a question-answering task. For evaluation, we use DeepEval by Confident AI [8], an open-source framework designed for LLM evaluation. DeepEval leverages LLMs and other NLP models to measure performance. In our study, GPT-3.5-turbo generates answers, with evaluation metrics focusing on contextual accuracy and relevance in both retrieval and generation stages. The following formulas are taken from DeepEval for evaluation. Additionally, we compare retrieval time across different strategies.

**Contextual Precision**

It measures how well relevant nodes are ranked higher in the retrieval context.

$$\text{CP} = \frac{1}{\text{Rel. Nodes}} \sum_{k=1}^{n} \left( \frac{\text{Rel. Nodes to } k}{k} \times r_k \right)$$

where $r_k$ is 1 for relevant nodes, 0 otherwise.

---

[5]LangChain
[6]Sentence Embedding: all-MiniLM-L6-v2
[7]Gemini Flash 1.5
[8]https://www.confident-ai.com

## Contextual Recall

The metric evaluates the ability of the system to capture relevant information:

$$CR = \frac{\text{Attributable Statements}}{\text{Total Statements}}$$

## Contextual Relevancy

It measures the overall relevance of the retrieval context with respect to the query:

$$CRel = \frac{\text{Relevant Statements}}{\text{Total Statements}}$$

## Answer Relevancy

Answer Relevancy evaluates the relevance of the generated output:

$$AR = \frac{\text{Relevant Statements}}{\text{Total Statements}}$$

## Faithfulness

Faithfulness measures how factually accurate the output is:

$$Faithfulness = \frac{\text{Truthful Claims}}{\text{Total Claims}}$$

## Retrieval Time

The Retrieval Time RT is defined as the total time taken to retrieve the context and generate the final answer for a query:

$$RT = t_{\text{end}} - t_{\text{start}}$$

These evaluation metrics allow us to compare the trade-offs between semantic integrity, retrieval effectiveness, and computational efficiency across different chunking approaches.

## 6 Results and Analysis

### 6.1 Results

Table 2 shows the chunk counts for different techniques. RSC achieves the best balance between granularity and coherence. In contrast, the Recursive Character Text Splitter generates the highest number of chunks due to its character-based splitting, while Semantic Chunking produces the fewest, resulting in larger segments. This balance reflects an important trade-off in RAG system design. Excessive chunking can inflate the retrieval space, leading to fragmented context. While larger chunks provide broader context, they increase the risk of irrelevant retrieval, hallucinations,

Table 2: Number of Chunks Formed by Each Chunking Method Across Datasets.

| Dataset | Recursive Character | Semantic | RSC (Proposed) |
|---|---|---|---|
| BBC News | 12,674 | 3,844 | 8,115 |
| SQuAD | 1,258 | 2,327 | 2,343 |
| QuAC | 2,464 | 2,307 | 4,121 |
| NewsMatrix-71 | 3,793 | 2,961 | 5,474 |

and longer retrieval times. RSC finds a middle ground, ensuring semantic integrity while maintaining meaningful chunk sizes. By keeping the chunk count within an optimal range, RSC improves contextual relevancy, as further supported by the downstream performance metrics in Table 3.

Table 3 presents the comparative performance of chunking techniques on the question-answering task across multiple datasets. The proposed Recursive Semantic Chunking consistently outperforms other techniques, particularly in Contextual Relevancy and Total Score, while maintaining an optimal balance between chunk size and retrieval efficiency.

The performance of chunking techniques across the datasets reveals interesting trends as shown in Figure 2. The best results are observed in SQuAD and NewsMatrix-71. SQuAD, achieving the highest Total Score under RSC, highlights the advantage of semantically coherent segmentation in structured question-answering datasets. NewsMatrix-71 achieves the highest Contextual Relevancy with RSC, demonstrating its effectiveness in handling diverse and large-scale articles.

In contrast, QuAC performs the worst, particularly under Semantic Chunking and Recursive Semantic Chunking. This is likely due to its conversational nature, which demands deeper contextual understanding.

While RSC does not lead in Answer Relevancy across all datasets, it is an important metric for evaluating end-to-end RAG performance. It consistently achieves top performance in Total Score and Contextual Relevancy. It is important to note that Answer Relevancy may be influenced by factors beyond chunking quality, such as the formulation of user queries (Sclar et al., 2024) or reasoning behavior of the language model during generation (Jiang et al., 2025). In contrast, Contextual Relevancy more directly reflects the quality and alignment of retrieved content with the query, making it a

Table 3: Performance Metrics for Different Chunking Techniques Across Datasets. Scores are out of 50, except Total Score (out of 250). Retrieval time is measured in seconds.
**Abbreviations:** RC = Recursive Character, RSC = Recursive Semantic Chunking, Avg Retv Time = Average Retrieval Time
**Bold values indicate the highest performance for each metric.**

| Metric | RC | Sem | RSC (Proposed) |
|---|---|---|---|
| **BBC News** | | | |
| Answer Relevancy | **45.89** | 43.89 | 41.97 |
| Answer Faithfulness | 38.51 | 35.18 | **42.55** |
| Contextual Recall | 43 | 45.5 | **46.33** |
| Contextual Precision | 47.02 | 44.82 | **48.98** |
| Contextual Relevancy | 11.40 | 8.78 | **11.56** |
| Total Score ↑ | 185.83 | 178.17 | **191.39** |
| Avg Retv Time(s) ↓ | 0.721 | 0.799 | **0.716** |
| **NewsMatrix-71** | | | |
| Answer Relevancy | **47.92** | 47.81 | 47.08 |
| Answer Faithfulness | **43.96** | 43.93 | 40.11 |
| Contextual Recall | 46.33 | **46.5** | 45.67 |
| Contextual Precision | 48.5 | 47.26 | **48.83** |
| Contextual Relevancy | 13.94 | 14.71 | **19.83** |
| Total Score ↑ | 200.65 | 200.21 | **201.52** |
| Avg Retv Time(s) ↓ | 0.72 | **0.71** | **0.71** |
| **SQuAD** | | | |
| Answer Relevancy | 47.28 | 46.67 | **48.59** |
| Answer Faithfulness | 44.98 | 43.71 | **46.5** |
| Contextual Recall | **50** | 49 | **50** |
| Contextual Precision | 47.09 | 47.99 | **47.99** |
| Contextual Relevancy | 17.7 | 20.09 | **20.12** |
| Total Score ↑ | 207.05 | 207.46 | **213.2** |
| Avg Retv Time(s) ↓ | 0.97 | 0.97 | **0.96** |
| **QuAC** | | | |
| Answer Relevancy | **45.4** | 44.69 | 43.67 |
| Answer Faithfulness | 41.675 | **44.25** | 43.63 |
| Contextual Recall | 47.08 | 45.33 | **48.58** |
| Contextual Precision | 43.67 | 45.16 | **45.76** |
| Contextual Relevancy | **12.47** | 9.64 | 9.38 |
| Total Score ↑ | 190.29 | 189.07 | **191.01** |
| Avg Retv Time(s) ↓ | **0.62** | 0.65 | 0.64 |

stronger indicator of chunking effectiveness.

Overall, among the chunking techniques, RSC achieves the highest Total Score across all datasets. The recursive breakdown mechanism in RSC ensures that large chunks do not negatively impact RAG tasks. Additionally, Contextual Relevancy improves significantly with RSC, as evident in datasets like BBC News (11.56) and NewsMatrix-

71 (19.83), demonstrating its capability to maintain semantic coherence.

These findings suggest the impact of the type and structure of data on the chunking techniques. However, in comparison, RSC is the most effective among the baseline chunking techniques.



Figure 2: Performance Comparison of Chunking Techniques Across Datasets

## 6.2 Analysis

To evaluate the impact of Recursive Semantic Chunking on retrieval efficiency and chunk coherence, we conduct performance analysis across multiple datasets. The evaluation uses datasets of varying structures such as structured question-answering datasets (SQuAD, QuAC) and unstructured large-scale datasets (BBC News, NewsMatrix-71). It ensures that our findings are generalizable across multiple RAG tasks.

**Study on Propositional Segmentation**

We conduct a study to analyze the effect of propositional segmentation incorporated in our proposed chunking technique. The hypothesis is that propositional segmentation enhances Contextual Relevancy.

To validate our hypothesis, we experiment by including propositional segmentation in RSC and compare the results. For this case study, we employ the BBC News dataset. The comparison of results is presented in Table 4

The results confirm that propositional segmentation improves Contextual Relevancy (11.56 to 16.09). However, it is to be that improvement comes at the cost of increased retrieval time (from 0.716s to 0.8183s). In addition, it also has a computational overhead to convert all the sentences into propositions before they can be passed on for

Table 4: Comparison of RSC with and without Propositional Segmentation on BBC News Dataset.

| Metric | RSC Without Propositions | RSC With Propositions |
|---|---|---|
| Answer Relevancy | **41.97** | 42.95 |
| Answer Faithfulness | **42.55** | 39.51 |
| Contextual Recall | **46.33** | 45.14 |
| Contextual Precision | **48.98** | 47.99 |
| Contextual Relevancy | 11.56 | **16.09** |
| Total Score ↑ | 191.39 | **191.68** |
| Avg Retv Time(s) ↓ | **0.716** | 0.8183 |

chunking. However, it is an interesting area of study for the future.

**Challenges of Agentic Chunking**

Although not included as a formal baseline, we initially explored Agentic Chunking to assess the viability of LLM-based chunking pipelines. However, due to its high computational demand, it is excluded from comparative evaluation. Details of Agentic Chunking are mentioned in Section 2. Since the Agentic approach operates at the propositional level, so for this technique, on average, each proposition requires 6 to 7 calls to the LLM for chunk assignment and metadata updates. To start with, we use this technique on the BBC dataset. The dataset contained more than 75,000 propositions, but after 8 hours of processing, only 1,500 propositions were successfully assigned to chunks. Due to the high computational overhead, we discontinued the experimentation. Hence, high computational cost makes this approach impractical for large-scale datasets.

Despite its inefficiencies, Agentic Chunking may become viable in the future as LLMs improve in speed and affordability. However, for now, RSC provides a far more efficient and scalable solution.

The results and analysis confirm that RSC enhances retrieval efficiency and semantic coherence. Additionally, our findings highlight a new direction with propositional segmentation, which improves Contextual Relevancy. Overall, RSC consistently outperforms both Recursive Character Text Splitter and Semantic Chunking in Total Score and Contextual Relevancy, making it the preferred approach for RAG generation pipeline. Moving forward, future work will focus on optimizing propositional segmentation to reduce retrieval time, ensuring that the benefits of enhanced semantic coherence do not come at the expense of computational overhead.

## 7 Conclusion

Our work offers a targeted contribution to optimizing the chunking process in RAG-based systems. The proposed technique, Recursive Semantic Chunking maintains a balance between retrieval efficiency and context relevancy. The novelty of RSC lies in the recursive nature of the proposed method dynamically adjusting the chunk size and going beyond the traditional approaches. The results, evaluated against the traditional techniques i.e. recursive character split, semantic and agentic techniques highlight the superiority of the proposed methodology. Additionally, its robustness is validated across structured question-answering datasets and unstructured large-scale datasets, with evaluation based on relevancy, retrieval quality, and time efficiency. The evaluation is based on relevancy, retrieval quality and time efficiency. These findings have significant implications for RAG-based applications such as medical, finance, legal, and education etc. Looking forward, the retrieval time will be further optimized with respect to Recursive Semantic Chunking on varied datasets.

## Limitation

The scope of this study is limited to textual data, and it can be widened to more complex document types which may include tables, codes etc. In addition, Recursive Semantic which depends on propositions provides a new direction. However, its high computational cost, despite yielding improved results, highlights the need for a more efficient and scalable approach.

## References

LangChain AI. 2024. Langchain. GitHub repository.

Mohammad Alkhalaf, Ping Yu, Mengyang Yin, and Chao Deng. 2024. Applying generative ai with retrieval augmented generation to summarize and extract key clinical information from electronic health records. *Journal of Biomedical Informatics*, 156:104662.

Tong Chen, Hongwei Wang, Sihao Chen, Wenhao Yu, Kaixin Ma, Xinran Zhao, Hongming Zhang, and Dong Yu. 2024. Dense X retrieval: What retrieval granularity should we use? In *Proceedings of the 2024 Conference on Empirical Methods in Natural Language Processing*, pages 15159–15177, Miami, Florida, USA. Association for Computational Linguistics.

Eunsol Choi, He He, Mohit Iyyer, Mark Yatskar, Wentau Yih, Yejin Choi, Percy Liang, and Luke Zettlemoyer. 2018. QuAC: Question answering in context. In *Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing*, pages 2174–2184, Brussels, Belgium. Association for Computational Linguistics.

Matthijs Douze, Alexandr Guzhva, Chengqi Deng, Jeff Johnson, Gergely Szilvasy, Pierre-Emmanuel Mazaré, Maria Lomeli, Lucas Hosseini, and Hervé Jégou. 2024. The faiss library. *Preprint*, arXiv:2401.08281.

Ruitao Feng, Xudong Hong, Mayank Jobanputra, Mattes Warning, and Vera Demberg. 2024. Retrieval-augmented modular prompt tuning for low-resource data-to-text generation. In *Proceedings of the 2024 Joint International Conference on Computational Linguistics, Language Resources and Evaluation (LREC-COLING 2024)*, pages 14053–14062, Torino, Italia. ELRA and ICCL.

FullStackRetrieval. 2024. Agentic chunker. Accessed: 2024-09-14.

Derek Greene and Pádraig Cunningham. 2006. Practical solutions to the problem of diagonal dominance in kernel document clustering. In *Proceedings of the 23rd International Conference on Machine Learning*, ICML '06, page 377–384, New York, NY, USA. Association for Computing Machinery.

Shiming He, Yu Hong, Shuai Yang, Jianmin Yao, and Guodong Zhou. 2024. Demonstration retrieval-augmented generative event argument extraction. In *Proceedings of the 2024 Joint International Conference on Computational Linguistics, Language Resources and Evaluation (LREC-COLING 2024)*, pages 4617–4625, Torino, Italia. ELRA and ICCL.

Yi Jiang, Sendong Zhao, Jianbo Li, Haochun Wang, and Bing Qin. 2025. GainRAG: Preference alignment in retrieval-augmented generation through gain signal synthesis. In *Proceedings of the 63rd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 10746–10757, Vienna, Austria. Association for Computational Linguistics.

LangChain. 2023. Recursively split by character. Accessed: 2024-09-14.

LangChain. 2024. Langchain documentation. Accessed: 2024-10-31.

LangChain. 2024. Semantic chunker. Accessed: 2024-09-14.

Patrick Lewis, Ethan Perez, Aleksandra Piktus, Fabio Petroni, Vladimir Karpukhin, Naman Goyal, Heinrich Küttler, Mike Lewis, Wen-tau Yih, Tim Rocktäschel, Sebastian Riedel, and Douwe Kiela. 2020.

Retrieval-augmented generation for knowledge-intensive nlp tasks. In *Proceedings of the 34th International Conference on Neural Information Processing Systems*, NIPS '20, Red Hook, NY, USA. Curran Associates Inc.

Jerry Liu. 2022. LlamaIndex.

Puneet Mathur, Zhe Liu, Ke Li, Yingyi Ma, Gil Karen, Zeeshan Ahmed, Dinesh Manocha, and Xuedong Zhang. 2024. DOC-RAG: ASR language model personalization with domain-distributed co-occurrence retrieval augmentation. In *Proceedings of the 2024 Joint International Conference on Computational Linguistics, Language Resources and Evaluation (LREC-COLING 2024)*, pages 5132–5139, Torino, Italia. ELRA and ICCL.

Carlo Merola and Jaspinder Singh. 2025. Reconstructing context: Evaluating advanced chunking strategies for retrieval-augmented generation. *Preprint*, arXiv:2504.19754.

Pranav Rajpurkar, Jian Zhang, Konstantin Lopyrev, and Percy Liang. 2016. SQuAD: 100,000+ questions for machine comprehension of text. In *Proceedings of the 2016 Conference on Empirical Methods in Natural Language Processing*, pages 2383–2392, Austin, Texas. Association for Computational Linguistics.

WA Sahlman, AM Ciechanover, and E Grandjean. 2023. Khanmigo: Revolutionizing learning with genai. *Harvard Business School Case*, pages 824–059.

Melanie Sclar, Yejin Choi, Yulia Tsvetkov, and Alane Suhr. 2024. Quantifying language models' sensitivity to spurious features in prompt design or: How i learned to start worrying about prompt formatting. *Preprint*, arXiv:2310.11324.

Spurthi Setty, Katherine Jijo, Eden Chung, and Natan Vidra. 2024. Improving retrieval for rag based question answering models on financial documents.

L. Siddharth and Jianxi Luo. 2024. Retrieval augmented generation using engineering design knowledge. *Knowledge-Based Systems*, 303:112410.

R. Teja. 2023. Evaluating the ideal chunk size for a rag system using llamaindex. Accessed: 2024-09-14.

Antonio Jimeno Yepes, Yao You, Jan Milczek, Sebastian Laverde, and Renyu Li. 2024. Financial report chunking for effective retrieval augmented generation. *Preprint*, arXiv:2402.05131.

# Towards Robust Urdu Aspect-based Sentiment Analysis through Weakly-Supervised Annotation Framework

**Zoya Maqsood[1]**    **Seemab Latif[1]**    **Rabia Latif[2]**

[1]School of Electrical Engineering and Computer Science,
National University of Sciences and Technology (NUST), Islamabad, Pakistan
{zmaqsood.phdcs17seecs, seemab.latif}@seecs.edu.pk
[2]College of Computer and Information Sciences (CCIS),
Prince Sultan University, Riyadh, Saudi Arabia
rlatif@psu.edu.sa

## Abstract

Aspect-Based Sentiment Analysis (ABSA) remains largely unexplored in low-resource languages like Urdu due to the absence of large-scale, publicly available, and domain-diverse annotated corpora. Additional challenges like the scarcity of lexical resources, unstructured Urdu websites, and linguistic complexities, further hinder corpus development. These limitations create a critical bottleneck that prevents robust Urdu ABSA systems from being deployed in practical scenarios. We address this gap by proposing a weakly supervised framework that automates corpus annotation (∼10K Budget tweets) leveraging seed-based pattern matching with dynamic window analysis. Through a comparative analysis of Large Language Models (LLMs), and human annotations on expertly curated datasets, we further demonstrate the inherent complexity of Urdu ABSA. Suboptimal results from a conventional LSTM model that achieved a mean performance of 0.52 precision, 0.49 recall, and 0.50 F1 score across various ABSA tasks validate this challenge. In short, this work establishes a scalable and cost-effective annotation framework that advances ABSA research for Urdu and similar low-resource languages.

## 1 Introduction

Aspect-Based Sentiment Analysis (ABSA) is a fine-grained Opinion Mining (OM) domain that evaluates sentiment toward specific attributes of entities, offering valuable insights for customer feedback analysis, product benchmarking, and market trend monitoring (Zhang et al., 2022). ABSA comprises four key elements: aspect category *(c)*, aspect term *(a)*, opinion term *(o)*, and sentiment polarity *(p)*. Figure 1 illustrates these elements through an annotation example of a customer's review. Examining single or multiple combinations of these elements to understand opinions in diverse scenarios gives rise to various ABSA tasks (Maqsood, 2023). For



Figure 1: Aspect Based Sentiment Analysis

instance, the extraction of aspect categories constitutes the Aspect Category Detection (ACD) task, whereas sentiment analysis over these categories leads to the Aspect Category Sentiment (ACS) task. Similarly, evaluating sentiment toward explicit aspect terms constitutes the Aspect Sentiment Classification (ASC) task.

A fundamental prerequisite for any OM system is an accessible benchmark corpus of annotated reviews (Zhou et al., 2019; Hu et al., 2021). This requirement becomes particularly acute for low-resource languages like Urdu, where despite substantial social media presence, available ABSA datasets remain inadequate characterized by non-public availability, absence of benchmark standards, sparse annotations, and limited-domain coverage (Rani and Anwar, 2020; Ahmad and Wan, 2021). While manual annotation of ABSA elements becomes prohibitively expensive for large-scale datasets containing multi-aspect sentences in various domains. Additional challenges include the

scarcity of lexical resources, prevalent use of non-standard encoding in Urdu web content, unique linguistic features, and informal language on social media (Khattak et al., 2021). These constraints collectively impede corpus development, creating significant barriers in building robust Urdu ABSA models and complicating the adaptation of existing methodologies (Zhou et al., 2019; Liu et al., 2020; Zhang et al., 2022).

Besides, leveraging weak supervision has demonstrated potential in the realm of social media mining (Maqsood, 2023; Tekumalla and Banda, 2023). Although weak labels may not achieve manual-level precision but they enable rapid dataset expansion and robust model training especially when combined with a subset of high-quality manual labels (Zhang et al., 2022). Despite the success of Large Language Models (LLMs) like GPT-4.0 and DeepSeek in capturing linguistic patterns, these approaches have not been widely explored in existing literature, particularly for dataset annotations in Urdu. While, applying English-centric models to translated Urdu tweets exacerbates the issue, yielding poor results due to translation quality limitations (Zhang et al., 2021).

This work pioneers Urdu ABSA by introducing a weakly supervised annotation framework that automates labeling of all core ABSA elements for the 'Budget' domain, overcoming dataset scarcity without costly manual effort. Our systematic evaluation reveals LLMs (GPT-4, DeepSeek) limited transferability to Urdu, while experiments demonstrate our method's superiority over them. To our knowledge, this constitutes the first comprehensive study of such techniques for Urdu. Baseline LSTM experiments further highlight Urdu-specific ABSA challenges, underscoring the need for advanced architectures. Our key contribution addresses Urdu's critical resource gap through scalable dataset creation methodology that eliminates manual annotation bottleneck to facilitate fine-grained Urdu ABSA.

## 2 Related Work

This section discusses the existing Urdu datasets developed for opinion mining tasks, analyzing their annotation methodologies, and domain applicability.

### 2.1 Opinion Mining Datasets in Urdu

Researchers contributed to the field of Urdu sentiment analysis by presenting annotated corpora, but most focus on document- or sentence-level sentiment classification rather than fine-grained ABSA. Early efforts, such as those by Rani and Anwar (2020), introduced a manually annotated corpus of 10,000 tweets from sports domains (cricket and football), labeling aspects, categories, and polarities. However, the absence of opinion term annotations limits applications of ABSA tasks. Similarly, ul Haq et al. (2020) presented a corpus of 8,760 political tweets with polarity and four category labels but did not annotate aspect terms or opinion expressions, restricting deeper sentiment analysis. Moreover, their dataset is not publicly available and labeled manually, hindering scalability and reproducibility.

To address the scarcity of ABSA-specific resources, Ahmad and Wan (2021) translated the SemEval-2014 ABSA dataset (2951 restaurant and 4721 laptop reviews) into Urdu, providing aspect terms, polarities, and category labels. While this enables some ABSA experimentation, the reliance on machine translation raises concerns about linguistic accuracy and cultural relevance. Other datasets, such as Ghafoor et al. (2023) introduced SentiUrdu1M dataset (1 million tweets), leveraging large-scale emoticon-based labeling but remain unsuitable for ABSA due to their document-level granularity. Similarly, Amjad et al. (2021) curated a dataset of 3,564 tweets for threat detection, but its binary classification focus makes it irrelevant for aspect-level sentiment tasks. Beyond Twitter data, researchers have collected Urdu reviews from blogs and news platforms. (Mukhtar and Khan, 2018; Mukhtar et al., 2017; Khan et al., 2021; Rehman and Bajwa, 2016) developed datasets with manually annotated sentiment labels at document-level and suffer from limited domain coverage (e.g., movies, electronics). Additionally, many of these datasets are not publicly available, and their annotation methodologies are often poorly documented, reducing their utility for ABSA research.

In short, existing Urdu sentiment analysis datasets lack fine-grained annotations, suffer from small sizes and narrow domains, and use inconsistent annotation methodologies. Most rely on either biased translations or labor-intensive manual labeling, which impedes scalability. Furthermore, existing resources neglect weakly-supervised

approaches, while available multilingual models and LLMs remain under-evaluated. Overall, these limitations underscore the dire need for comprehensively annotated Urdu ABSA datasets in several domains by combining both manual and automated annotation methods. This hybrid methodology ensures both high-quality annotations and efficient scalability, ultimately enabling advanced techniques for progress of ABSA in Urdu language.

## 3 Dataset

We collected approximately 13,000 tweets related to Pakistan's budgetary domain between May and July 2020 using Twitter's Standard API. Due to API constraints, tweets were gathered in daily batches, limited to a 7-day historical window, with a maximum of 100 tweets per query and 180 requests per 15-minute interval. The search queries focused on trending budgetary discourse in Pakistan, incorporating hashtags such as '#Budget2020', '#PakistanEconomy', and '#Commerce'. The dataset provides a comprehensive representation of public sentiments and economic debates surrounding Pakistan's budget during the unprecedented COVID-19 lockdown period.

### 3.1 Pre-processing

The collected tweets underwent an extensive three-stage pre-processing pipeline to ensure data quality and linguistic consistency.

**Tweet Level:** We performed Unicode normalization to address Arabic script variations, removed punctuations, and social media artifacts (emojis, hashtags, URLs, mentions) using regular expressions. We eliminated duplicate entries and truncated excessive consecutive repetitions (e.g., reducing "سلیکٹڈ بحث..." (selected budget selected budget...) to "سلیکٹڈ بحث") (selected budget) to maintain textual conciseness.

**Token Level:** After conducting a systematic comparison of tokenization approaches Qi et al. (2020), Ali (2020), Vasiliev (2020) and space-based methods, we preferred UrduHack for its superior performance on informal Urdu text. Use of informal language and noise on social media limit the effectiveness of language-specific tokenizers, introducing abnormal tokens. We analyzed incorrect tokens to identify the inherent patterns of their abnormalities and normalized them accordingly. This includes splitting merged stopwords (e.g., تھیاک ←اک تھی), reducing character repetitions in

misspelled words (e.g., پاکسسستان ←پاکستان), and eliminating word repetitions (e.g., فریفریفری ←فری).

**Character Level:** The final processing step validated individual characters against Urdu Unicode ranges and removed residual artifacts (e.g., cleaning "***" and normalizing "هوـنـڈ کاc002u" to "هوـنـ کا").

This hierarchical pre-processing approach, documented comprehensively in Zoya et al. (2023), resulted in dataset of approximately 10,000 tweets.

### 3.2 Dataset Variants

We created three versions of the dataset, introducing variation in the annotation process, as outlined below:

**Bronze Standard Dataset (BS):** This dataset is a raw output without manual curation from our weakly supervised annotation system.

**Silver Standard Dataset (SS):** This represents a refined version of the 'BS' dataset. The corpus underwent a meticulous validation process combining automated consistency checks with expert human verification to ensure higher annotation quality. This approach filtered out erroneous labels generated by our weakly supervised methods and resulted in an 13% reduction of the original dataset labels.

**Gold Standard Dataset (GS):** The GS dataset was constructed through rigorous manual annotation by three native Urdu speakers with expertise in NLP. From the SS corpus, we selected a representative subset of 3000 tweets for fine-grained annotations. Three annotators followed strict annotation guidelines of Pontiki et al. (2014) standards, with only labels receiving consensus from at least two annotators being retained. The GS corpus serves as a reliable ground truth for evaluating model performance on Urdu ABSA tasks, while also revealing additional linguistic patterns not captured in the initial SS annotations. The statistics about these datasets have been described in Table 1.

| Dataset | Tweets | Asp_Cat. | Asp_Terms | Opinion_Terms |
|---------|--------|----------|-----------|---------------|
| Bronze  | 9693   | 14       | 5179      | 5456          |
| Silver  | 8949   | 14       | 4247      | 5364          |
| Gold    | 3000   | 14       | 1126      | 1410          |

Table 1: Statistics of Datasets with Distinct Values

## 4 Methodology

We present our methodology for annotating Urdu datasets for ABSA. First, we highlight the limita-

tions of LLMs for this task, followed by our custom framework designed to address these challenges.

## 4.1 LLMs Limitations for Dataset Annotation

The utilization of the GPT 4.0 and DeepSeek models for dataset annotation in Urdu revealed several challenges. Firstly, the model encountered challenges in thoroughly capturing all aspects and sentiment words present in tweets. Secondly, an inherent instability in labeling responses was observed, as the model exhibited varying results for the same query when executed multiple times. Thirdly, the issue of selecting irrelevant words alongside sentiment and aspect terms introduced a lack of uniformity, necessitating post-processing efforts for pruning. Fourthly, the model tended to repeat sentiment terms within aspect terms or vice versa. Fifthly, breaking down tweets into shorter chunks did not significantly improve their response quality. Sixthly, the model demonstrated a tendency to ignore rare words and occasionally overlook crucial aspects. In conclusion, LLMs exhibited limitations in fully grasping the context. A representative case of tweet annotation generated by the LLM in Figure 4 (see appendix).

## 4.2 Dataset Annotation Framework for ABSA in Urdu

Our dataset labeling approach encompasses two fundamental phases: ACD and the annotation of Aspect-Opinion-Sentiment (AOS) triplet. Initially, ACD was completed through topic modeling and clustering techniques. Subsequently, the identification of triplet components within tweets was carried out through methods like pattern mining and a bidirectional window-based labeling strategy.

### 4.2.1 Aspect Category Detection

We used pre-trained sentence transformers Reimers and Gurevych (2019) to generate embeddings and applied both Top2Vec (Angelov, 2020) and traditional clustering algorithms (Ackermann et al., 2014; Frey and Dueck, 2007) to identify nuanced subtopics. Our analysis revealed optimal cluster counts (39) based on cosine similarity metrics and cluster validation techniques (Kaoungku et al., 2018; Yuan and Yang, 2019). Notably, Top2Vec initially predicted 54 topics, but these were ultimately clustered within the same range.

To reduce cluster overlap, we performed graph-based analysis, where edges represented cosine similarities between embeddings. Edge weights were set to 0 for similarities below a threshold of 0.7, ensuring that only highly similar tweets were grouped together while preserving distinct terms across clusters. However, some topics (clusters) exhibited irrelevance like synonymous terms or polysemy of less substantial words. Such problematic topics featuring highly coherent terms could form distinct clusters, leading to favourable scores in standard metrics. Conversely, some significant topics might be overlooked due to lower coherence or similarity scores between words, particularly if such topics cover diverse perspectives not covered well in the coherence metrics reference corpus. To address these limitations, we incorporated a manual curation step to refine and consolidate topics. From the generated clusters, we selected distinct categories (Table 3 in section 8) and further subdivided broad topics by analyzing top topic words (Table 4 in section 8). For example, 'Social Welfare' was divided into 'Education', 'Agriculture', 'Health', and 'Social Programs'.

### 4.2.2 Aspect-Opinion-Sentiment Triplet Annotation

This triplet annotation process is divided into four fundamental stages, which include word classification, seed enrichment, tweet labeling, and evaluation, as discussed below:

**Words Classification:** The selected topics consist of a mixture of terms related to aspect and opinion that require further classification. To systematically categorize these terms (seed terms), we adopted a straightforward yet effective approach: nouns were designated as aspect, while adjectives were treated as indicators of opinion. The selection process for seed terms emphasized domain relevance, frequency, and diversity, ensuring the chosen nouns were explicit and closely related to core topics. To enhance relevance, we excluded irrelevant or ambiguous terms, rare occurrences, verbs, adverbs, and any expressions introducing sentiment bias or lacking clear aspect association. This rigorous selection process resulted in a refined lexicon of aspect and opinion seeds, with the complete workflow detailed in Figure 2.

**Seeds Enrichment:** Given the limited coverage of initial seed words for annotating all tweets within each subtopic, we employed multiple strategies to expand and refine our seed term collection. As illustrated in Figure 3, our enrichment approach incorporated sentiment lexicons, active learning,

Figure 2: Aspect Categories and Words Classification Process

pattern mining, and embedding-based methods. **Sentiment Lexicon:** We used an existing Urdu



Figure 3: Sentiment Classification From Predicted Topics

sentiment lexicon [42] to identify polar expressions within our tweet corpus. This process yielded 1,322 positive and 1,395 negative terms that overlapped between the lexicon and our budget-related tweets. However, the lexicon exhibited notable limitations: its coverage of domain-specific fiscal terminology was incomplete, and its formal vocabulary often mismatched the informal expressions and morphological variations prevalent in social media. Similar challenges emerged when we attempted to use translated lexicons intended for sentiment analysis in the English language.

**Multiform words:** We took into account various word forms in our seed terms, including singular and plural forms, such as "قیمت" (price) in singular and "قیمتوں" or "قیمتیں" (prices) in plural. Urdu's rich inflectional system, where words vary by tense, gender, number, and loanword integration, renders conventional lemmatization and stemming ineffective. For instance, contextual variants (e.g., verb conjugations or gendered forms) lack reliable root-mapping rules. Moreover, Urdu has a diverse vocabulary with numerous loanwords and context-dependent variations that further complicate such tasks. Consequently, we excluded this step to preserve semantic precision given the absence of robust Urdu-specific linguistic tools.

**Active Learning:.** We utilized an active learning approach and created a preliminary dataset consisting of hundred short tweets (5-10 words each). Subsequently, we conducted manual labelling with a specific focus on AOS triplets. This step provided valuable insights for various types of words beyond the seed words and their contextual relationships within the domain. Additionally, we noted the prevalence of multi-word phrases as opposed to single words for seed terms. We quantified phrase frequencies and their sentiment associations, iteratively expanding the seed lexicon to include high-impact multi-word terms. This comprehensive examination not only enriched our seed inventory but also our understanding of the multifaceted language used in the dataset.

**Pattern Mining:** Based upon the manually labeled data from active learning, we identified recurring patterns that encompassed consecutive domain-specific words and Urdu case markers (کا (ka), کے (kay), کی (ki), کو (ko), میں (mein), پر (par), سے (se), نے (nay)). We developed a hybrid pattern mining approach combining rule-based and statistical techniques. This integrated approach revealed important multi-word expressions that served as more precise indicators of aspects and opinions compared to conventional single-word seeds. We first analyzed recurring syntactic structures involving domain-specific terms paired with Urdu case markers. Matching these patterns against tweets is depicted in Algorithm 1 and a comprehensive list of extracted patterns is provided in Table 5 (Section 8). We then implemented a sequential pattern mining algorithm with minimum support thresholds to discover statistically significant co-occurring word sequences, prioritizing longer phrases that captured more nuanced meanings. The extracted patterns enabled us to automatically identify aspect-opinion pairs in new tweets. For instance, in the structure "[X] ka [Y]", X was classified as the aspect term and Y as the opinion term. Sentiment polarity was then assigned to these newly discovered opinion terms through contextual analysis, leading to the formation of AOS triplets (detailed in Appendix Algorithm 2).

**Embeddings:** We utilized the pre-trained embedding model FastText to identify the top 10 words that exhibited the highest cosine similarity with our seed terms, particularly focusing on expanding our set of opinion words. Additionally, we considered terms returned by the Top2vec model that exhib-

**Algorithm 1** Patterns Matching Algorithm

---

1: **procedure** EXTRACTPATTERNS(budget_tweets)
2:    pattern ← (\w+\sکے بجٹ\s)
3:    extracted_patterns ← EmptyList()
4:    **for** each tweet in budget_tweets["tweet_text"] **do**
5:
6:       matches ← FindAllMatches(pattern, tweet)
7:       **for** each match in matches **do**
8:          pattern_text ← match.group(1)
9:
10:          AddPatternToExtractionList(extracted_patterns, pattern_text)
11:       **end for**
12:    **end for**
13:    **return** extracted_patterns
14: **end procedure**

---

ited similarity with seed terms by surpassing the 0.5 threshold in similarity score. We selectively kept words that fell within the categories of aspect or opinion-related terms. Any words failing to meet these criteria were excluded from further consideration. Exemplary instances have been presented in the Table 6 (Appendix).

**Labeling Tweets:** Initially, we annotated using mined patterns with analogous structures, resulting in the creation of AOS triplets while accounting for sentiment reversals caused by negators (e.g., no, not). Annotated example can be seen in Figure 5 (in Appendix), proved effective for contiguous word patterns but limited for non-adjacent term relationships. To address this limitation, we introduced a window-based annotation strategy consisting of two main steps: seeds cartesian product with sentiment polarity assignment and seeds co-occurrence analysis.

**Seeds Cartesian Product with Sentiment Polarity Assignment:** In this phase, we performed a Cartesian product operation between the aspect seeds and sentiment seeds to form their pairs (a, o). Despite that sentiment polarity was already predefined in lexicons for numerous opinion words, several pairs underwent cross-validation due to domain-specific variations or informal language use. As the sentiment of the same opinion word may vary based on its association with different aspect words. For example, the term increment is considered positive when associated with salary but negative when linked with poverty.

**Window-based Strategy:** We implemented a dynamic window-based approach to detect co-occurring aspect-opinion (a, o) pairs within tweets. The tweet segmentation process entailed setting a token threshold of 15 words from both the be-

ginning and end. Subsequently, we systematically examined the co-existence of (a, o) pairs within these segments. Meanwhile, we addressed negators when they co-occurred within a segment in the context of the (a, o) pair, and selectively inverted the sentiment for that particular (a, o) pair occurrence in the given tweet. Likewise, we advanced to the next segment by adjusting the window after every five words. Eventually, we reconstructed the original tweet and gathered all unique AOS triplets from every segment of a tweet. This comprehensive methodology allowed us to discern nuanced sentiment variations associated with (a, o) pairs within the dynamic context of tweets. The entire process is summarized below and labeled tweet result is presented in Figure 6 (Appendix).

$$segment\_length = 15$$

$$window\_size = 5$$

$$S_i = \text{Segment}(t\_i, segment\_length)$$

$$C_{i,j} = \text{CoExistence}(S_{i,j}, (a, o))$$

$$N_{i,j} = \text{NegatorHandling}(C_{i,j}, \text{negators})$$

$$S_i = \text{NextSegment}(S_i, window\_size)$$

$$R_i = \text{Reconstruct}(S_i, N_i)$$

$$AOS_i = \text{GatherTriplets}(R_i)$$

**Evaluation:** To assess label quality, we employed dual evaluation approaches as given below:
**Automated Evaluation:** We compared our methodology's outputs against both LLMs annotations and Gold-standard (GS) labels on identical tweet samples. Our analysis revealed that our proposed approach mostly outperformed LLMs in AOS triplet accuracy.

**Human Evaluation:** We developed a weakly-supervised validation protocol addressing multi-aspect tweets where window-based strategies occasionally produced spurious aspect-opinion associations in case of multiple aspects. The validation process involved: (1) categorizing tweets by presence of aspect complexity (single/multiple), (2) cross-referencing novel multi-aspect pairs with pre-labeled single-aspect examples and pattern-mined results, and (3) manual verification of unmatched pairs on a sample representing at least 2% of the tweets containing each such pair. If more than 50% of the labels were deemed accurate in the chosen sample, we retained them as final labels.

Finally, to ensure label consistency across datasets, we performed comparative analysis by identifying tweet overlaps between all dataset variants and discrepancies were compared against the GS labels. Then the F1 measure and accuracy were computed (Table 2) as defined in (Pontiki et al., 2014) and expressed below:

$$F1 = \frac{2 \cdot P \cdot R}{P + R} \qquad (1)$$

where precision (P) and recall (R) were determined as:

$$P = \frac{|SS \cap GS|}{|SS|} \qquad (2)$$

$$R = \frac{|SS \cap GS|}{|GS|} \qquad (3)$$

$$Acc. = \frac{|GS \cap SS|}{|SS \cup GS|} \qquad (4)$$

| Label | Acc. | P | R | F1 |
|---|---|---|---|---|
| Aspect | 69.9 | 91.7 | 74.6 | 82.3 |
| Opinion | 71.4 | 89 | 80.9 | 84.8 |
| Polarity | 73.6 | 88.5 | 74.4 | 80.8 |
| Category | 86.3 | 93.1 | 89.8 | 91.4 |

Table 2: Scores of evaluation measures on annotated dataset labels.

## 5 Experimental Set-Up

### 5.1 Tasks

We performed experiments on three key ABSA tasks, as given below:

**Aspect Category Detection (ACD):** Identifying the categories for each tweet from a set of predefined aspect categories.

**Aspect Category Sentiment (ACS):** Sentiment polarity classification (positive/negative/neutral) toward detected aspect categories.

**Aspect Sentiment Classification (ASC):** Sentiment polarity analysis targeting explicit aspect terms.

### 5.2 Model

We implemented LSTM as our baseline model initialized with 300-dimensional FastText embeddings. The model was trained with a batch size of 32, hidden state dimension of 300, and the adam optimizer (learning rate = 0.001) for 100 epochs. To ensure robustness, we ran five training repetitions

using categorical cross-entropy loss. Hyperparameters were tuned via Grid search, testing epochs [10, 50, 100, 300], embedding dimensions [100, 300], learning rates [0.001, 0.01, 0.0001], batch sizes [16, 32, 64, 128], and dropout rates [0.2, 0.3, 0.5], with early stopping (patience = 5) and stratified 5-fold cross-validation. The hyperparameter grid values are chosen based on optimal LSTM performance observed in sentiment analysis-related studies (Kumar et al., 2021; Naqvi et al., 2021).

### 5.3 Dataset Distribution

We implemented a rigorous train-test split (Table 7 in Section 8) on the SS dataset to maintain proportional representation of both aspect categories and sentiment polarities. The partitioning preserved identical distributions of positive, negative, and neutral sentiment labels across training (75%) and testing (25%) subsets for each aspect category. The equal percentage distribution provides a balanced representation for classifier training and fosters robust model development by minimizing biases through learning from comparable instances across various aspect categories.

### 5.4 Results Analysis

The LSTM baseline results reveal a consistent performance trend across datasets (results in Appendix 8). For ACD, the GS achieves strong performance at 100 epochs, while SS and BS show gradual improvements, peaking at 0.596 and 0.562 accuracy, respectively. This aligns with the high F1 scores (91.4 for Category, 82.3 for Aspect) in Table 2, confirming that our annotation framework produces usable labels. In ACS, the GS reached near-ceiling macro-F1 (0.877) by 50 epochs, whereas SS and BS plateau at ~0.490.59 F1. This reflects the challenge of sentiment polarity prediction. The SS consistent lead over BS dataset justifies our refinement step, though both trail Gold due to inherent noise. For ASC, all datasets struggle (F1 < 0.35), mirroring the difficulty of fine-grained sentiment analysis. The marginal gains with more epochs suggest the LSTMs limited capacity to resolve ambiguities. Traditional LSTM is viable for coarse tasks (ACD) but face limitations in sentiment-related tasks. However, the LSTM model was intentionally selected as a lower-bound baseline to assess the discriminative strength of annotation quality and task difficulty, without the confounding influence of pretraining or large-scale parameters in advanced architectures. Despite balanced splits,

macro-F1 scores highlight challenges from label imbalance, multi-label learning, and Urdu's morphological complexity. Progressive performance gains from (Bronze→Silver→Gold) highlight annotation quality as a stronger factor than model complexity.

## 6 Discussion

The proposed weakly-supervised framework demonstrates significant advancements in Urdu ABSA by overcoming the critical bottleneck of manual annotation in dataset creation. The multidimensional annotation requirements, encompassing all ABSA elements, render fully manual annotation impractical for scalable model development due to its time-consuming nature and human labor requirements. Our framework automates this process, starting with a seed-based approach for high-precision in noisy, code-mixed Urdu social media text and mitigate limitation of domain coverage through iterative enrichment using lexicon expansion, syntactic patterns, and contextual embedding strategies. This dynamic refinement transforms static seeds into a robust, domain-adaptive seeds inventory suited for low-resource and informal text settings. Thus, the core strength lies in the novel integration of context-aware seed expansion and morphologically-sensitive preprocessing, which collectively reduce annotation costs.

Furthermore, the method demonstrates robust capability in handling Urdu's linguistic complexities through its hybrid approach combining n-gram pattern matching with dynamic window labeling. This approach effectively identifies multi-word aspects, such as "پٹرول کی قیمت" (petrol price), and successfully resolves polarity inversion cases by incorporating negation scope detection. Additionally, an automated validation pipeline was introduced that minimize human effort to maintain label quality. The limited variation with Gold-Standard dataset underscores the significance of high-quality annotations from our proposed method. Comparative analysis with prevailing LLMs reveals the proposed framework achieves substantially better performance for annotation task in Urdu, especially for aspect and opinion terms extraction tasks. These advancements establish a practical foundation for Urdu ABSA where fully supervised approaches remain infeasible due to resource constraints. Regarding classification results, the performance of conventional models like LSTM across

ABSA tasks and datasets are emphasized. Despite extended training, the baseline LSTM's limited improvement reveals its inability to capture Urdu's linguistic nuances in ABSA tasks. There were instances where additional epochs do not yield significant gains, suggesting a potential saturation point in the models learning curve. Although our evaluation is constrained to the budget domain due to the availability of gold-standard annotations, framework's core components such as linguistic and syntactic pattern rules, clustering mechanism, and seed augmentation are domain-independent and easily adaptable to other domains.

### 6.1 Limitations

The proposed dataset annotation methodology endeavors to address many challenges, yet certain issues persist. Sentiments occurring beyond segment window lengths are occasionally overlooked, although this is mitigated by considering segments from multiple positions within tweets. The exclusion of tweets lacking seed terms may inadvertently dismiss relevant sentiment expressions. Overlapping labels or spurious associations may emerge occasionally when a sentiment word applies in multiple perspectives or simultaneously relates to multiple aspects within a tweet. In cases like sarcasm, where the same word is employed in diverse contexts (positively or neutrally), priority is determined based on its frequency of occurrence. Polysemous terms labeling (e.g., for budget pass vs. approach) risks errors, highlighting needs for context-aware rules.

## 7 Conclusion

Our research introduced a novel weak supervision methodology for creating a benchmark dataset in Urdu ABSA. We shed light on the inherent challenges in ABSA for under-resourced languages and made a significant contribution to addressing the resource scarcity in Urdu ABSA. Our dataset encompasses tweets within the budget domain and was annotated at four distinct levels: aspect, opinion, sentiment, and category levels. The consistently high F1 scores across all label annotations demonstrate the proposed method's effectiveness in producing high-quality. Through a detailed comparative analysis involving LLMs and human annotations based on expertly curated datasets, we illuminated the intricate nature of our proposed dataset. Empirical evaluations utilizing LSTM model showed limita-

tions of conventional methods for various ABSA subtasks and laid the groundwork for future advancements in ABSA techniques for Urdu.

# 8 Future Work

We aim to generalize our methodology to expand ABSA dataset annotations into other domains. Our focus will extend to advanced deep learning techniques, moving beyond basic LSTM models for diverse ABSA tasks in Urdu. We plan to conduct fine-tuning pre-trained models on an extended dataset across various domains for a comprehensive understanding of Urdu sentiment expressions. In summary, our future trajectory involves leveraging advanced techniques, annotating diverse datasets, and refining models for domain-specific applications, ultimately enhancing Urdu ABSA tools.

# References

Marcel R Ackermann, Johannes Blömer, Daniel Kuntze, and Christian Sohler. 2014. Analysis of agglomerative clustering. *Algorithmica*, 69:184–215.

Naveed Ahmad and Jing Wan. 2021. Aspect based sentiment analysis for urdu. In *2021 6th International Conference on Computational Intelligence and Applications (ICCIA)*, pages 309–313.

Ikram Ali. 2020. Urduhack: A python library for Urdu language processing. *CoRR*, abs/2010.06810. ArXiv: 2010.06810.

Maaz Amjad, Noman Ashraf, Alisa Zhila, Grigori Sidorov, Arkaitz Zubiaga, and Alexander Gelbukh. 2021. Threatening language detection and target identification in Urdu tweets. *IEEE Access*, 9:128302–128313.

Dimo Angelov. 2020. Top2vec: Distributed representations of topics. *CORR*, abs/2008.09470. ArXiv: 2008.09470.

Brendan J Frey and Delbert Dueck. 2007. Clustering by passing messages between data points. *Science*, 315(5814):972–976.

Abdul Ghafoor, Ali Shariq Imran, Sher Muhammad Daudpota, Zenun Kastrati, Sarang Shaikh, and Rakhi Batra. 2023. Sentiurdu-1m: A large-scale tweet dataset for urdu text sentiment analysis using weakly supervised learning. *PLOS ONE*, 18(8):e0290779.

Mengting Hu, Shiwan Zhao, Honglei Guo, Chao Xue, Hang Gao, Tiegang Gao, Renhong Cheng, and Zhong Su. 2021. Multi-label few-shot learning for aspect category detection. In *Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing (Volume 1:*

*Long Papers)*, pages 6330–6340, Online. Association for Computational Linguistics.

Nuntawut Kaoungku, Keerachart Suksut, Ratiporn Chanklan, Kittisak Kerdprasop, and Nittaya Kerdprasop. 2018. The silhouette width criterion for clustering and association mining to select image features. *International journal of machine learning and computing*, 8(1):69–73.

Lal Khan, Ammar Amjad, Noman Ashraf, Hsien-Tsung Chang, and Alexander Gelbukh. 2021. Urdu sentiment analysis with deep learning methods. *IEEE Access*, 9:97803–97812.

Asad Khattak, Muhammad Zubair Asghar, Anam Saeed, Ibrahim A Hameed, Syed Asif Hassan, and Shakeel Ahmad. 2021. A survey on sentiment analysis in urdu: A resource-poor language. *Egyptian Informatics Journal*, 22(1):53–74.

Avinash Kumar, Aditya Srikanth Veerubhotla, Vishnu Teja Narapareddy, Vamshi Aruru, Lalita Bhanu Murthy Neti, and Aruna Malapati. 2021. Aspect term extraction for opinion mining using a hierarchical self-attention network. *Neurocomputing*, 465:195–204.

Haoyue Liu, Ishani Chatterjee, MengChu Zhou, Xiaoyu Sean Lu, and Abdullah Abusorrah. 2020. Aspect-based sentiment analysis: A survey of deep learning methods. *IEEE Transactions on Computational Social Systems*, 7(6):1358–1375.

Zoya Maqsood. 2023. Weakly supervised learning for aspect based sentiment analysis of urdu tweets. In *Proceedings of the 8th Student Research Workshop associated with the International Conference Recent Advances in Natural Language Processing*, pages 78–86.

Neelam Mukhtar and Mohammad Abid Khan. 2018. Urdu sentiment analysis using supervised machine learning approach. *International Journal of Pattern Recognition and Artificial Intelligence*, 32(02):1851001.

Neelam Mukhtar, Mohammad Abid Khan, and Nadia Chiragh. 2017. Effective use of evaluation measures for the validation of best classifier in urdu sentiment analysis. *Cognitive Computation*, 9:446–456.

Uzma Naqvi, Abdul Majid, and Syed Ali Abbas. 2021. Utsa: Urdu text sentiment analysis using deep learning methods. *IEEE Access*, 9:114085–114094.

Maria Pontiki, Dimitris Galanis, John Pavlopoulos, Harris Papageorgiou, Ion Androutsopoulos, and Suresh Manandhar. 2014. SemEval-2014 task 4: Aspect based sentiment analysis. In *Proceedings of the 8th International Workshop on Semantic Evaluation (SemEval 2014)*, pages 27–35, Dublin, Ireland. Association for Computational Linguistics.

Peng Qi, Yuhao Zhang, Yuhui Zhang, Jason Bolton, and Christopher D Manning. 2020. Stanza: A python natural language processing toolkit for many human languages. *CoRR*, abs/2003.07082. ArXiv: 2003.07082.

Sadaf Rani and Waqas Anwar. 2020. Resource creation and evaluation of aspect based sentiment analysis in urdu. In *Proceedings of the 1st Conference of the Asia-Pacific Chapter of the Association for Computational Linguistics and the 10th International Joint Conference on Natural Language Processing: Student Research Workshop*, pages 79–84.

Zia Ul Rehman and Imran Sarwar Bajwa. 2016. Lexicon-based sentiment analysis for urdu language. In *2016 Sixth International Conference on Innovative Computing Technology (INTECH)*, pages 497–501.

Nils Reimers and Iryna Gurevych. 2019. Sentence-BERT: Sentence embeddings using Siamese BERT-networks. In *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)*, pages 3982–3992, Hong Kong, China. Association for Computational Linguistics.

Ramya Tekumalla and Juan M Banda. 2023. Using weak supervision to generate training datasets from social media data: a proof of concept to identify drug mentions. *Neural Computing and Applications*, 35(25):18161–18169.

Ehsan ul Haq, Sahar Rauf, Sarmad Hussain, and Kashif Javed. 2020. Corpus of aspect-based sentiment for urdu political data. In *Mexican International Conference on Artificial Intelligence*, pages 37–40. Springer.

Yuli Vasiliev. 2020. *Natural Language Processing with Python and SpaCy: A Practical Introduction*. No Starch Press.

Chunhui Yuan and Haitao Yang. 2019. Research on k-value selection method of k-means clustering algorithm. *J*, 2(2):226–235.

Jieyu Zhang, Cheng-Yu Hsieh, Yue Yu, Chao Zhang, and Alexander Ratner. 2022. A survey on programmatic weak supervision. *CoRR*, abs/2202.05433. ArXiv: 2202.05433.

Wenxuan Zhang, Ruidan He, Haiyun Peng, Lidong Bing, and Wai Lam. 2021. Cross-lingual aspect-based sentiment analysis with aspect term code-switching. In *Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing*, pages 9220–9230.

Jie Zhou, Jimmy Xiangji Huang, Qin Chen, Qinmin Vivian Hu, Tingting Wang, and Liang He. 2019. Deep learning for aspect-level sentiment classification: survey, vision, and challenges. *IEEE access*, 7:78454–78483.

Zoya, Seemab Latif, Rabia Latif, Hammad Majeed, and Nor Shahida Mohd Jamail. 2023. Assessing urdu language processing tools via statistical and outlier detection methods on urdu tweets. *ACM Transactions on Asian and Low-Resource Language Information Processing*, 22(10):1–31.

# Appendix

**Prompt:** *Extract all aspect expressions with their corresponding opinion expressions and sentiment polarity (positive/ negative/neutral) in given text:*

## Aspect-Opinion-Sentiment Analysis

**Text:**

"ملک معاشی طور پر دیوالیہ ہو چکا ہے اور بجٹ بنانے میں دقت پیش آ رہی ہے۔"

**Extracted Pairs:**

| Aspect Term (پہلو) | Opinion Term (رائے) | Sentiment Polarity (جذباتی) (کیفیت) |
|---|---|---|
| ملک کی معاشی حالت (Economic condition) | دیوالیہ ہو چکا ہے (Bankrupt) | Strong Negative (انتہائی) (منفی) |
| بجٹ سازی (Budget preparation) | دقت پیش آ رہی ہے (Facing difficulty) | Negative (منفی) |

Figure 4: DeepSeek response to Annotate Tweet

| S# | Label | Top 15 Words |
|---|---|---|
| 1 | سماجی بہروگرام Social Welfare | ہیلتھ, روزگار, وائرس, کسان, وباء, زراعت, اصلاحات, ایجوکیشن, تعلیم, ڈاکٹر, فنکاروں, صحت, طلبہ, شعبہ, محکمہ 'Health', 'Employment', 'Virus', 'Farmer', 'Epidemic', 'Agriculture', 'Reforms', 'Education', 'Learning', 'Doctor', 'Artists', 'Health', 'Students', 'Department' |
| 2 | معیشت Economy | معاشی, مالی, مہنگا, پالیسی, سود, سالانہ, بحران, خسارے, معیشت, ڈالر, سود, معاشیات, اقتصادی, ترقیاتی, خزانہ 'Economic', 'Financial', 'Expensive', 'Policy', 'Interest', 'Annual', 'Crisis', 'Losses', 'Economy', 'Dollar', 'Interest', 'Economics', 'Economic', 'Development', 'Treasury' |
| 3 | میڈیا Media | تفصیلات, اطلاعات, پریس, رپورٹ, معلومات, میڈیا, صحافی, نیوز, خبر, احتجاج, تقریر, اعلان, عوامی, اداروں, حکومتیں 'Details', 'Information', 'Press', 'Report', 'Knowledge', 'Media', 'Journalist', 'News', 'Protest', 'speech','Report', 'Announcement', 'Public', 'Institutions', 'Governments' |
| 4 | سیاست Politics | سرکاری, وفاقی, صدر, حکمران, اپوزیشن, پاکستان, وزراء, سیاسی, ملک, ریاست, ایوان, اسمبلی, کابینہ, سیاست, پارٹی 'Government', 'Federal', 'President', 'Rulers', 'Opposition', 'Pakistan', 'Ministers', 'Political', 'Country', 'State', 'Assembly', 'Cabinet', 'Politics', 'Party' |
| 5 | مذهب Religion | اسلام, محمد, اسلامی, مسلم, گلوکار, اللہ, شریف, شیخ, علما, مدارس, دین, بجٹ, حکومت, سود, مدینہ 'Islam', 'Muhammad', 'Islamic', 'Muslim', 'Singer', 'Allah', 'Sharif', 'Scholar', 'Scholars', 'Schools', 'Religion', 'Budget', 'Government', 'Interest', 'Medina' |
| 6 | دفاع Defense | بم, پنشن, خطرہ, جنگ, افواج, پاک, دفاع, دشمن, دہشت, پولیس, ٹیکنالوجی, اصلاحات, اداروں, مراعات, عالمی 'Bomb', 'Pension', 'Threat', 'War', 'Forces', 'Pakistan', 'Defense', 'Enemy', 'Terror', 'Police', 'Technology', 'Reforms', 'Institutions', 'Consideration', 'Global' |
| 7 | منصوبہ Project | سبسڈی, پروجیکٹ, منصوبے, مختص, فنڈ, اخراجات, روپے, پیسہ, خرچ, قرضوں, کرپشن, فری, معاشرے, ترقیاتی, پالیسیوں 'Subsidy', 'Project', 'Projects', 'Specialized', 'Fund', 'Expenditure', 'Rupees', 'Money', 'Expense', 'Loans', 'Corruption', 'Free', 'Society', 'Development', 'Policy' |

Table 3: Selected subtopics derived from predicted clusters and topic modeling

156

| S# | Aspect Category | Attributes |
|---|---|---|
| 1 | معیشت<br>Economy | بجٹ، عمومی، اخراجات، مہنگائی، قیمت، محصول، قرض، سود<br>budget, general, expenditure, inflation, price, revenue, debt, interest |
| 2 | وفاق<br>Federal | بجٹ، جنرل، وزیراعظم، صدر، کابینہ، اسمبلی، حکومت<br>budget, general, prime minister, president, cabinet, assembly, government |
| 3 | تعلیم<br>Education | بجٹ، جنرل، ادارے، استاد، طالب علم<br>budget, general, institutions, teacher, student |
| 4 | صحت<br>Health | بجٹ، جنرل، ڈاکٹر، صحت، ہسپتال، دیکھ بھال، وبا<br>budget, general, doctor, health, hospital, care, epidemic |
| 5 | زراعت<br>Agriculture | بجٹ، جنرل، زراعت، کسان، فصل، ٹڈی<br>budget, general, Agriculture, farmer, crop, locust |
| 6 | سماجی بہروگرام<br>Social Welfare Program | بجٹ، جنرل، اصلاحات، بینظیر انکم ٹیکس ترقیاتی پروگرام، اصلاحاتی پروگرام، انکم ٹیکس سپورٹ پروگرام<br>budget, general, reforms, Benazir income tax, development program,<br>reform program, income tax support program |
| 7 | دفاع<br>Defense | بجٹ، جنرل، فوجی، تحفظ، حملہ<br>budget, general, military, protection, attack |
| 8 | مذہب<br>Religion | بجٹ، جنرل، مذہب، مومن، علماء، مقدس مقامات<br>budget, general, religion, believers, scholars, holy places |
| 9 | سیاسی جماعت<br>Political Party | بجٹ، جنرل، پارٹی، پالیسی، کانفرنس، اپوزیشن<br>budget, general, party, policy, conference, opposition |
| 10 | قیادت<br>Leadership | بجٹ، جنرل، لیڈر، چیئرمین، کرپشن<br>budget, general, leader, chairman, corruption |
| 11 | صوبائی<br>Provincial | بجٹ، جنرل، صوبہ (پنجاب، سندھ، بلوچستان، پختونخوان) حکومت، کابینہ، اسمبلی<br>budget, general, provinces (Punjab, Sindh, Balochistan, Pakhunkhawan)<br>govt., cabinet, assembly |
| 12 | عوام<br>Public Dynamics | بجٹ، جنرل، امیر، غریب، روزگار، تنخواہ، پنشن<br>budget, general, rich, poor, employment, salary, pension |
| 13 | میڈیا<br>Media | بجٹ، جنرل، صحافی، میڈیا، خبریں، چینل، رپورٹ، آرٹسٹ<br>budget, general, journalist, media, news, channel, report, artist |
| 14 | جنرل<br>Miscellaneous | بجٹ<br>budget |

Table 4: Conclusive sub-categories of budget topic

| Patterns | 'کا بجٹ'<br>(Budget of) | 'دوست بجٹ'<br>(Friend's Budget) | 'دشمن بجٹ'<br>(Enemy's Budget) | 'میں کمی'<br>(Decrease in) |
|---|---|---|---|---|
| Phrases | امراء کا بجٹ<br>(Budget of Aristocrats)<br>تباھی کا بجٹ<br>(Budget of Destruction)<br>خسارے کا بجٹ<br>(Budget of Loss)<br>مافیا کا بجٹ<br>(Mafia's Budget)<br>تعلیم دوست بجٹ<br>(Education-Friendly Budget) | انسان دوست بجٹ<br>(Human-Friendly Budget)<br>عوام دوست بجٹ<br>(Public-Friendly Budget)<br>تعلیم دوست بجٹ<br>(Education-Friendly Budget)<br>غریب دوست بجٹ<br>(Poor-Friendly Budget)<br>نوجوان دوست بجٹ<br>(Youth-Friendly Budget) | انسانیت دشمن بجٹ<br>(Inhumane Budget)<br>برآمدات دشمن بجٹ<br>(Incomes-Enemy Budget)<br>صحت دشمن بجٹ<br>(Health-Enemy Budget)<br>مزدور دشمن بجٹ<br>(Labor-Enemy Budget)<br>مسلم دشمن بجٹ<br>(Muslim-Enemy Budget) | حکومتی اخراجات میں کمی<br>(Reduction in Government Expenditure)<br>تعلیمی اخراجات میں کمی<br>(Reduction in Educational Expenditure)<br>بجٹ خسارہ میں کمی<br>(Reduction in Budget Loss)<br>فیسوں میں کمی<br>(Reduction in Fee)<br>پٹرولیم قیمتوں میں کمی<br>(Reduction in in Petroleum Prices) |
| Total | 240 | 20 | 39 | 59 |

Table 5: Phrases extracted by the Pattern Mining

**Algorithm 2** : MineSequentialPatterns
___

1: **procedure** MINESEQUENTIALPATTERNS(budget_tweets)
2:     stopwords ← LoadStopwords()                                              ▷ Load stopwords
3:     ps ← PrefixSpanAlgo(data)                                 ▷ Initialize pattern mining algorithm
4:     min_support ← 20                                                    ▷ Set minimum support
5:                                                 ▷ Mine frequent patterns with minimum support
6:     result ← ps.Frequent(min_support)
7:                                                                            ▷ Filter patterns
8:     filtered_patterns ← FILTERPATTERNS(result, stopwords)
9:                                                                ▷ Display and store patterns
10:    obt_patterns ← DISPLAYANDSTOREPATTERNS(filtered_patterns)
11:    **return** obt_patterns                                      ▷ Return the obtained patterns
12: **end procedure**

___

1: **function** FILTERPATTERNS(result, stopwords)
2:     filtered_patterns ← []                                           ▷ List for filtered patterns
3:     **for** each (support, pattern) in result **do**
4:                                                                   ▷ Check if pattern is valid
5:         **if** ISPATTERNVALID(pattern, stopwords) **then**
6:                                                                 ▷ Keep valid pattern to list
7:             filtered_patterns.append((pattern, support))
8:         **end if**
9:     **end for**
10:    **return** filtered_patterns                                  ▷ Return the filtered patterns
11: **end function**

___

1: **function** ISPATTERNVALID(pattern, stopwords)
2:                                                                 ▷ Check length of pattern
3:     **if** Length(pattern) > 1 **then**
4:                                                               ▷ Count stopwords in pattern
5:         stopwords_count ← COUNTSTOPWORDS(pattern, stopwords)
6:         **if** stopwords_count ≤ 1 **then**
7:             is_subpattern ← False                             ▷ Initialize flag for subpattern
8:             **for** each (_, other_pattern) in result **do**
9:                                                     ▷ Check if pattern is subset of other pattern
10:                **if** pattern ≠ other_pattern & ISSUBSET(pattern, other_pattern) **then**
11:                    is_subpattern ← True
12:                    **break**                                 ▷ Exit loop if subpattern is found
13:                **end if**
14:            **end for**
15:            **if** not is_subpattern **then**
16:                **return** True
17:            **else**
18:                **return** False
19:            **end if**
20:        **else**
21:            **return** False
22:        **end if**
23:    **else**
24:        **return** False
25:    **end if**
26: **end function**

| Seeds | Top 10 similar words |
|-------|---------------------|
| قرض (Loan) | 'قرضہ', 'قرضوں', 'کایس', 'قرضدار', 'ادهار', 'مقروض', 'پُرقرض', 'اُدهار', 'Risky' 'Loan', 'Loans', 'Debts', 'Cabinet', 'Debtor', 'Interest', 'Indebted', 'Owing', 'Debt', 'Risky' |
| حکومت (Govt.) | 'حکومتوں', 'حکومتی', 'نوازحکومت', 'حکومتیں', 'بشارحکومت', 'کوصوبے', 'ہیںحکومت', 'ہےحکومت', 'وزارت' 'Governments', 'Governmental', 'Nawaz Govt.', 'Governments', 'Bashar Govt.', 'In Govt.', 'In Govt.', 'Are in Govt.', 'Is in Govt.', 'Ministry' |
| بجٹ (Budget) | '0084ارب', 'روڈنئ', 'شیڈوبجٹ', 'کابجٹ', 'کےریلیف', '05ارب', '04ارب', '74 کهرب57ارب', '005ارب', '006ارب' '4800 Billion', 'Rodney', 'Shadow Budget', 'Cabinet', 'Relief', '50 Billion', '40 Billion', '47 Billion 75 Million', '500 Billion', '600 Billion' |
| مہنگائی (Inflation) | 'پرمہنگائی', 'کومہنگائی', 'اورمہنگائی', 'پهرمہنگائی', 'مہنگائ', 'مہنگائی', 'هوشرباء', 'قیمتیں', 'قیمتوں', 'اورغریت' 'Hyperinflation', 'And Inflation', 'And Inflation', 'Then Inflation', 'Inflation', 'Inflation', 'Hosharba', 'Prices', 'Prices', 'And Poverty' |
| تنخواہ (Salary) | 'تنخوا', 'تنخواہیں', 'سےتنخواہ', 'تنخواہ', 'تنخواہیں', 'تنخواہ', 'تنخواہیں', 'تنخواہوں', 'تنخوائں', 'تنخواہو' 'Salary', 'Salaries', 'From Salary', 'Salary', 'Salaries', 'Salaries', 'Salary', 'Salaries', 'Salary', 'Salaries' |

Table 6: Most Similar words by FastText model

"Urdu Tweet": "غریب دشمن بجٹ نامنظور دهاندلی کی حکومت نامنظور پی ٹی ایم ایف بجٹ نامنظور مزدور دشمن بجٹ نامنظور عوام دشمن بجٹ نامنظور"

"Translated Tweet": "Anti-Poor budget disapproved. The government of rigging disapproved. PTIMF budget disapproved. Anti-labor budget disapproved. Anti-public budget disapproved.

```
"entries": [                                  {
    {                                             "Aspect": "بجٹ" (budget),
        "Aspect": "بجٹ" (budget),                 " Opinion": " غریب دشمن " (anti-public),
        "Opinion": "عوام دشمن" (anti-public),     "Category": "بجٹ" ( budget),
        "Category": "بجٹ" (budget),               "Polarity": "Negative"
        "Polarity": "Negative"                 },
    },                                         {
    {                                             "Aspect": " حکومت " (government),
        "Aspect": "بجٹ" (budget),                 " Opinion": " دهاندلی " (rigged),
        "Opinion": "نامنظور"(disapprove)",        "Category": " وفاق " (Federal),
        "Category": "بجٹ" (budget),               "Polarity": "Negative"
        "Polarity": "Negative"                 },
    },                                         {
    {                                             "Aspect": " حکومت " (government),
        "Aspect": "بجٹ" (budget),                 " Opinion": "نامنظور" (disapprove),
        "Opinion": " مزدور دشمن " (anti-labor),   "Category": " وفاق " (Federal),
        "Category": "بجٹ" (budget),               "Polarity": "Negative"
        "Polarity": "Negative"                 }
    },                                         ]
```

Figure 5: Pattern Mining: Aspect-Sentiment Labels Division Based on Identified Phrases-Similar color shows single pattern

"Urdu Tweet": " آئی ایم ایف کی ہدایت پر کٹھ پتلی وزیر اعظم عمران نیازی کا بجٹ پاکستان کے غریبوں کے خلاف
اعلان جنگ ہے لعنت ہے ایسی تبدیلی پر کامریڈ افتخار"

Translated Tweet: The puppet Prime Minister Imran Niazi's budget on the instructions of IMF is a declaration of war against the poor of Pakistan. Curse such a change, Comrade Iftikhar.

```
"entries": [
    {
        "aspect": "وزیر اعظم عمران نیازی کا بجٹ" (Budget of Prime Minister Imran Niazi),
        "sentiment": "غریبوں کے خلاف" (against the poor),
        "category": "بجٹ" (Budget),
        "polarity": "Negative"
    },
    {
        "aspect": "وزیر اعظم عمران نیازی کا بجٹ" (Budget of Prime Minister Imran Niazi),
        "sentiment": "لعنت" (Curse),
        "category": "بجٹ" (Budget),
        "polarity": "Negative"
    }
]
```

Figure 6: Unique Labels obtained by bidirectional Window-based Strategy

| Aspect Categories | Train | | | | | Test | | | | |
|---|---|---|---|---|---|---|---|---|---|---|
| | Positive | Negative | Neutral | Total | P(%) | Positive | Negative | Neutral | Total | P(%) |
| Education | 42 | 78 | 20 | 139 | 10% | 9 | 12 | 2 | 47 | 10% |
| Agriculture | 5 | 10 | 3 | 18 | 1% | 1 | 1 | 0 | 6 | 1% |
| Federal | 246 | 459 | 114 | 820 | 57% | 52 | 68 | 14 | 274 | 57% |
| Media | 2 | 4 | 1 | 7 | 0% | 1 | 2 | 0 | 3 | 1% |
| Economy | 13 | 23 | 6 | 42 | 3% | 3 | 3 | 1 | 14 | 3% |
| Provincial | 18 | 34 | 8 | 60 | 4% | 4 | 5 | 1 | 21 | 4% |
| Political party | 11 | 21 | 5 | 37 | 3% | 2 | 3 | 1 | 13 | 3% |
| Health | 38 | 70 | 18 | 125 | 9% | 8 | 10 | 2 | 42 | 9% |
| Project | 4 | 7 | 2 | 13 | 1% | 1 | 1 | 0 | 5 | 1% |
| Social Welfare Programs | 20 | 36 | 9 | 65 | 5% | 4 | 5 | 1 | 22 | 5% |
| Leadership | 9 | 16 | 4 | 29 | 2% | 2 | 3 | 0 | 10 | 2% |
| Defense | 19 | 35 | 9 | 63 | 4% | 4 | 5 | 1 | 21 | 4% |
| Religion | 2 | 4 | 1 | 8 | 1% | 1 | 1 | 0 | 3 | 1% |
| Miscellaneous | 2 | 3 | 1 | 4 | 0% | 0 | 1 | 0 | 2 | 0% |
| **Total** | 429 | 800 | 201 | 1430 | 100% | 91 | 120 | 24 | 483 | 100% |
| **Percentage(%)** | 30% | 56% | 14% | 75% | – | 29% | 59% | 12% | 25% | – |

Table 7: Polarity-Specific Aspects Categories Distribution in Train-Test Split

| Task | Epoch | Gold-Standard | | Bronze-Standard | | Silver-Standard | |
|---|---|---|---|---|---|---|---|
| | | Acc. | macro-F1 | Acc. | macro-F1 | Acc. | macro-F1 |
| Aspect Category Detection (ACD) | 10 | 0.660 | 0.681 | 0.437 | 0.417 | 0.449 | 0.439 |
| | 50 | 0.711 | 0.722 | 0.505 | 0.526 | 0.580 | 0.526 |
| | 100 | 0.732 | 0.733 | 0.528 | 0.506 | 0.596 | 0.545 |
| | 300 | 0.717 | 0.723 | 0.562 | 0.525 | 0.590 | 0.549 |
| Aspect Category Sentiment (ACS) | 10 | 0.669 | 0.802 | 0.531 | 0.456 | 0.524 | 0.536 |
| | 50 | 0.687 | 0.877 | 0.566 | 0.487 | 0.531 | 0.590 |
| | 100 | 0.706 | 0.877 | 0.571 | 0.493 | 0.556 | 0.571 |
| | 300 | 0.706 | 0.877 | 0.575 | 0.494 | 0.524 | 0.590 |
| Aspect Sentiment Classification (ASC) | 10 | 0.575 | 0.296 | 0.561 | 0.243 | 0.582 | 0.258 |
| | 50 | 0.577 | 0.304 | 0.577 | 0.304 | 0.583 | 0.281 |
| | 100 | 0.578 | 0.318 | 0.581 | 0.313 | 0.589 | 0.302 |
| | 300 | 0.583 | 0.321 | 0.583 | 0.318 | 0.591 | 0.318 |

Table 8: LSTM Average Results for five-runs on ABSA Tasks

# Efficient Continual Learning for Small Language Models with a Discrete Key-Value Bottleneck

**Andor Diera**
Ulm University
andor.diera@uni-ulm.de

**Lukas Galke**
University of Southern Denmark
galke@imada.sdu.dk

**Fabian Karl**
Ulm University
fabian.karl@uni-ulm.de

**Ansgar Scherp**
Ulm University
ansgar.scherp@uni-ulm.de

## Abstract

Continual learning remains a challenge across various natural language processing (NLP) tasks, as models updated with new training data often risk catastrophic forgetting of previously acquired knowledge. We introduce a discrete key-value bottleneck (DKVB) for encoder-only language models, enabling efficient continual learning through localized updates. Inspired by a discrete key-value bottleneck in vision, we consider new and NLP-specific challenges. We compare different bottleneck architectures for NLP and introduce a new, task-independent initialization technique for the discrete keys. We evaluate our DKVB for NLP in four continual learning scenarios and show that it alleviates catastrophic forgetting. Our experiments demonstrate that the proposed approach achieves competitive performance compared to popular continual learning methods while incurring lower computational costs. Furthermore, we show that DKVB remains effective even in challenging single-head continual learning scenarios where no task ID is provided.[1]

## 1 Introduction

Large language models are receiving increasing attention from the public due to their impressive zero-shot and few-shot abilities in a wide range of tasks (Brown et al., 2020). Yet, for easier tasks where there is enough training data for supervised fine-tuning, e. g., text classification, using smaller encoder-only language models is still preferable due to their often superior performance and lower computational requirements (Yuan et al., 2023; Yu et al., 2023; Qorib et al., 2024; Li et al., 2025). Compared to large general-purpose models, fine-tuned networks lack general portability to new conditions and have limited generalization beyond their training distribution (Luo et al., 2023). For many target applications in natural language processing (NLP), training and test data can have a

difference in the underlying distribution (Hupkes et al., 2023), and in the case of continual learning, the input distribution can change over time (Wang et al., 2024). To mitigate these challenges, different changes to model architectures and training regimens have been proposed (Biesialska et al., 2020; Ke and Liu, 2022; Wang et al., 2024). While many of these methods improve continual learning, they often require task-specific modules and computationally demanding extensions to the base model (Ke et al., 2021; Buzzega et al., 2020; Momeni et al., 2025).

In this work, we propose an adaptation of the Discrete Key-Value Bottleneck (DKVB) architecture (Träuble et al., 2023) to the field of NLP. Discretization techniques can improve generalization in neural networks without introducing new task-specific parameters, regularization functions, or memory buffers (Liu et al., 2021, 2023; Träuble et al., 2023). More specifically, the DKVB architecture has shown strong performance in low-resource, class incremental learning scenarios for computer vision. This is due to local, context-dependent updates on learnable discrete key-value pairs that prevent catastrophic forgetting in the models.

To address the challenges of adapting DKVB to NLP, we begin by analyzing how different variants of the discrete key-value bottleneck interact with pre-trained encoder-only language models in standard learning scenarios. In doing so, we tackle key challenges such as the high dimensionality of text representations, the choice of pooling strategies, and the design of an effective decoder head. Subsequently, we take the best-performing DKVB configurations and evaluate their performance in continual learning scenarios. Finally, we show that given a dictionary of discrete keys optimized on a general-purpose corpus, DKVB achieves similar effectiveness compared to leading continual learning approaches while requiring less training time. The main contributions of our paper are:

---

[1]Source code available at: github.com/drndr/dkvb_nlp

- We analyze different optimization techniques and architectures of a DKVB in NLP using BERT, RoBERTa, and DistillBERT.

- We compare our DKVB for NLP to baseline methods in continual learning scenarios, i. e., domain incremental, class incremental, and task-type incremental learning.

- We demonstrate that the DKVB alleviates catastrophic forgetting and is more efficient than most continual learning methods.

## 2 Related Work

### 2.1 Continual Learning

Sequentially learning multiple tasks remains a significant challenge in the field of deep learning. Standard neural networks trained on a new task tend to forget most of the knowledge tied to tasks they have previously learned, leading to the phenomenon commonly labeled as *catastrophic forgetting* (McCloskey and Cohen, 1989; Van de Ven and Tolias, 2019). On the other hand, leveraging knowledge learned from old tasks to improve performance on new tasks, known as *knowledge transfer*, is a highly sought-after capability in NLP (Ke and Liu, 2022). Since re-training a model from scratch is often expensive, various methods for continual learning have been proposed to handle these challenges. Existing approaches in continual learning can be categorized into five distinct families: regularization-based, optimization-based, replay-based, architecture-based, and instruction-based, with the latter being specific to large language models (Biesialska et al., 2020; Ke and Liu, 2022; Wang et al., 2024; Shi et al., 2024). A detailed description of these approaches can be found in Appendix A.

### 2.2 Discrete Representation Learning

Employing discrete variables in deep learning is challenging, as indicated by the prevalence of continuous latent variables in most research methods, even when the underlying modality inherently involves discrete elements (e. g., text data). Van Den Oord et al. (2017) were the first to show the viability of large-scale discrete neural representation learning through the use of vector quantization. Their Vector Quantized-Variational Autoencoder (VQ-VAE) model utilizes a discrete latent space and thus avoids the "posterior collapse" problem common in many VAE models when the decoder ignores the latent space of the encoder and relies

solely on the autoregressive properties of the input samples (Goyal et al., 2017). Subsequently, their methodologies have been widely employed in various applications, including audio (Borsos et al., 2023), videos (Yan et al., 2021), and anomaly detection (Marimont and Tarroni, 2021). More recently, discretization has been utilized for machine unlearning (Shah et al., 2023) and to improve disentangled representation learning (Noh et al., 2023) and robustness (Liu et al., 2021, 2023; Träuble et al., 2023). Discretization methods with bottlenecks have been shown to improve generalization in reinforcement learning (Liu et al., 2021, 2023), visual reasoning (Liu et al., 2023), and vision-based continual learning (Träuble et al., 2023).

## 3 A Discrete Key-Value Bottleneck for Encoder-only Language Models

The DKVB architecture as described in Träuble et al. (2023) is fundamentally model and task-agnostic, but so far has been only studied in the field of computer vision. The use of DKVB in language models poses new challenges, including the (i) sequential nature of the input data, the (ii) high dimensionality of the encoded representations, and the (iii) difference in commonly used pooling techniques between vision and language models. Below, we describe DKVB's base architecture, the key initialization process, and our proposed architectural adaptations and pre-experiments for finding the most suitable architectural variant.

### 3.1 Base Architecture and Key Initialization

The DKVB architecture follows three steps: encode input, process via a discrete bottleneck, and decode. Figure 1 show an overview of the architecture.
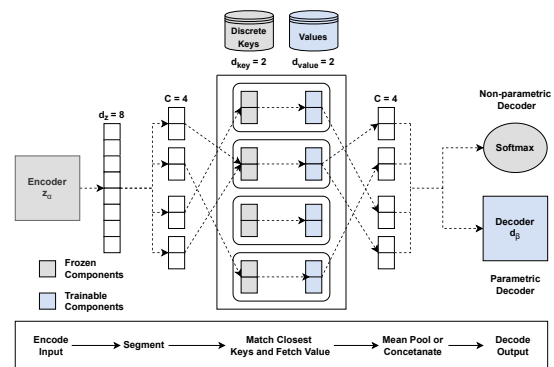


Figure 1: The base Discrete Key-Value Bottleneck.

In the first step, an encoder model projects input vector $x$ into a lower dimensional vector $z \in \mathbb{R}^{m_z}$. This is followed by pooling (if needed) and partitioning $z$ into $C$ separate heads of dimension $d_{key}$. Each head possesses a unique discrete key-value codebook of size $K$, where the keys are initialized before training and are mapped to randomly initialized trainable value codes. In the second step, each head is first quantized by fetching the closest key (based on L2 distance) from the corresponding head's codebook. Subsequently, the corresponding value code of dimension $d_{value}$ is retrieved for each head. Note the size of the bottleneck with respect to trainable parameter scales with the number of heads and codebook size. In the last step, the values are passed to a decoder to produce the final output. The decoder can be either parametric (with trainable weights) or non-parametric (by applying the softmax function to the mean pooled value codes).

The discrete keys of the bottleneck are initialized before training. Due to the 1-1 mapping between the keys and value codes, there is no gradient back-propagation between the values and keys. To ensure that the keys are broadly distributed in the feature space and have good representational power for given downstream tasks, they are first randomly initialized and then modified by using the encoded input samples as the basis for applying exponential moving average (EMA) updates (Van Den Oord et al., 2017). Alternatively, the keys can be initialized on input data different from the one in training, albeit with some decrease in downstream task performance (Träuble et al., 2023). After initialization, the keys are frozen and are not influenced by later changes in the input distribution shifts.

## 3.2 Architecture Adaptations for NLP

We introduce an adaptation of the DKVB architecture for the specific challenges in natural language processing. As argued above, these challenges are related to the high dimensionality of the data, pooling techniques, and decoding. We conduct pre-experiments with different architectures to find the most suitable bottleneck architecture variants and consider the following NLP-specific challenges:

**Dimensionality**   While natural language has an inherently discrete symbolic representation, text embeddings encode these discrete symbols into a continuous latent space (Muennighoff et al., 2023). This results in a high dimensional output $z \in \mathbb{R}^{t \times h}$, where $t$ is the token dimension (i.e., the number

of tokens in the fixed length input sequence) and $h$ is the hidden dimension. Previous experiments with DKVB were conducted on low dimensional image data that has been pooled before forwarding output $z \in \mathbb{R}^h$ to the bottleneck (Träuble et al., 2023). To address this difference, we design model variants with pooling applied before or after the bottleneck. Similarly, we experiment with creating the heads by partitioning hidden dimension $h$ and token dimension $t$ separately.

**Pooling Type**   Most modern convolutional networks in computer vision utilize max pooling as pooling operation (He et al., 2015). Max pooling retains the most important features in images but is less commonly used in NLP due to the loss of sequential information. The two most commonly used pooling techniques in NLP are mean pooling and pooling based on a special token (CLS). In mean pooling, the contextualized token embeddings are averaged out, while the CLS pooling utilizes a special token optimized to represent the whole sequence (Devlin et al., 2018). We include both variants in our architecture search.

**Decoding**   Decoders with adjustable weights offer more expressiveness than non-parametric decoders but are more sensitive to changes in the training conditions (Ostapenko et al., 2022). For simple tasks where linear mapping is sufficient, using just a softmax function as a non-parametric decoder might be appropriate. However, for many NLP tasks, it is crucial to capture complex patterns in the encoded representations (Wang et al., 2018). We include both approaches in our experiments. For the parametric decoder, we concatenate the value codes and feed them into a simple linear layer preceded by a dropout layer. In the non-parametric version, we apply mean pooling on the values and apply a softmax function on the pooled representation.

## 3.3 Analyzing DKVB Variants for NLP

We analyze different variants of the DKVB architecture in encoder-only language models. For the pre-experiments, we use two popular text classification datasets. The R8 dataset, which is a subset of the R21578 news dataset (Lewis, 1997) with 8 classes, and the Twenty Newsgroup (20ng) (Lang, 1995) which contains documents categorized into 20 newsgroups. We apply the standard train-test split for both datasets, as used in (Galke and Scherp,

Table 1: Accuracy and standard deviation (in subscript) of the different DKVB architecture variants on the R8 and 20ng datasets in a non-continual, standard learning setup, averaged over 5 runs.

| Decoder | Segmentation | Dataset: R8 | | | | Dataset: 20ng | | | |
| | | Pooling Before | | Pooling After | | Pooling Before | | Pooling After | |
| | | CLS | Mean | CLS | Mean | CLS | Mean | CLS | Mean |
|---|---|---|---|---|---|---|---|---|---|
| Parametric | hidden | $69.87_{1.07}$ | $91.54_{1.19}$ | $88.55_{1.05}$ | $\mathbf{96.04_{0.26}}$ | $19.26_{2.42}$ | $53.62_{1.00}$ | $48.35_{0.69}$ | $\mathbf{77.83_{0.89}}$ |
| | token | - | - | $88.06_{1.00}$ | $95.20_{1.21}$ | - | - | $44.65_{0.86}$ | $69.33_{0.96}$ |
| Non Parametric | hidden | $66.61_{0.29}$ | $92.18_{0.36}$ | $88.53_{0.22}$ | $94.24_{0.39}$ | $21.09_{1.31}$ | $55.93_{0.75}$ | $52.03_{0.24}$ | $73.51_{0.20}$ |
| | token | - | - | $64.39_{0.20}$ | $73.70_{0.20}$ | - | - | $10.95_{0.18}$ | $15.26_{1.28}$ |
| BERT (frozen) w/o DKVB | | $95.94_{0.18}$ | | | | $72.11_{0.49}$ | | | |
| BERT w/o DKVB | | $98.00_{0.34}$ | | | | $84.06_{0.53}$ | | | |

2022). We first use a frozen BERT model as the pre-trained encoder for DKVB and perform a hyper-parameter search on the number of epochs, batch size, and learning rates.

We report the performance of the best configurations. For learning rates, we found it is beneficial to have a high learning rate for the values layer. Additionally, in the case of the parametric decoder setup, a lower learning rate is applied to the decoder. We use a key dimension of $12$ and the number of key-value pairs of $4,096$ for the discrete bottleneck parameters as in (Träuble et al., 2023). Key initialization is done before training for three epochs with an EMA decay of 0.2. Alongside the different variants for the DKVB, we list the results of a fine-tuned BERT and a frozen BERT with a fine-tuned linear classifier on top for reference. This we consider as the upper bounds.

The test performance of the different architecture configurations can be seen in Table 1. The gap between the best-performing DKVB architecture and the fully fine-tuned BERT model is 2% on R8 and 7% on 20ng. The frozen BERT model achieved the same performance on R8 but attained 5% lower accuracy on 20ng compared to the best-performing DKVB variant. Overall the best performance was obtained by using the parametric decoder, applying mean pooling after the bottleneck, and using the hidden dimension as the base of the segmentation. To investigate the performance on other encoder-only language models, we experimented with RoBERTa (Liu et al., 2019) and DistilBERT (Sanh et al., 2019), and found the optimal bottleneck architecture to be the same (see Appendix C).

## 4 Continual Learning Settings

The goal of continual learning (CL) is to sequentially learn a function $f : X_k \rightarrow Y_k$ for all tasks

$k$ in sequence $K$. Each task $k$ has a training set $M_k = \{(x_i, y_i, d_i, t)\}_{i=1}^{N_k}$, where $x_i \in X_k$ is a training sample, $y_i \subseteq Y_k$ is a set of class labels, $d_i \subseteq D_k$ is the corresponding domain set (e. g., legal documents, movie reviews, news articles) of the sample, $t \in T_k$ is the task-type of the training set (e. g., sentiment analysis, topical classification, natural language inference etc), and $N_k$ is the number of samples in task $k$. To evaluate the DKVB architecture, we define three different incremental learning setups based on these components.

In the **Domain Incremental Learning (DIL)** setting, the task type and class labels are assumed to be consistent across all tasks. The domain of the input changes between tasks, with each task having a set of non-overlapping domains $D_k \cap D_{k'} = \emptyset$. A common DIL task-type in NLP is sentiment classification, where all tasks have the same class labels (i. e., positive, negative, neutral), but include samples from different source domains.

In the **Class Incremental Learning (CIL)** setting, each task has a set of non-overlapping classes $Y_k \cap Y_{k'} = \emptyset$. During testing, any previously learned class may be presented. CIL is generally considered the most challenging incremental learning scenario (Ke and Liu, 2022; Träuble et al., 2023). Apart from catastrophic forgetting and knowledge transfer, this setting includes the added complication of inter-task class separation (Kim et al., 2022). Inter-task class separation requires learning decision boundaries between the new task's classes and the classes from previous tasks without access to data from those previous tasks.

The main challenge in the **Task-type Incremental Learning (TIL)** setting lies in the varying task-types. While it is possible that the tasks also have non-overlapping input domains and class labels in these settings, what differentiates TIL from other

incremental learning scenarios are the disjoint task-types in each task $t_k \neq t_{k'}$. This task-type is not identical to the type of objective function used in training (i. e., classification loss or regression loss); rather, it defines the downstream task of the model, such as topical classification, sentiment analysis, or measuring semantic similarity. The scenario of using one model to learn different task-types has also been heavily researched in the field of multi-task learning (Crawshaw, 2020). The dominant approach in TIL is using a multi-head configuration with a separate head (or output layer) for each task. Since this decreases the probability of catastrophic forgetting, the main challenge in TIL is bi-directional knowledge transfer (Ke and Liu, 2022).

## 5 Experimental Setup

In our continual learning experiments, we compare the performance of DKVB to other CL methods in the three settings described above, namely TIL, DIL, and CIL. In addition, we adapt the challenging single-head CIL setup from Trauble et al. (Träuble et al., 2023) to topical text classification. We take the best-performing bottleneck architectures from the pre-experiments and apply the same bottleneck parameters and hyperparameters. We use accuracy as the primary evaluation metric in all our experiments and present the average performance and standard deviation over five runs. These runs involve random initialization and randomized task sequence order. Additionally, we report the average per epoch runtime of each method. Details about the implementation, hyperparameters, and bottleneck parameters can be found in the Appendix B.

**Datasets**  For the main experiments, we use three datasets, two of which have also been used in Ke et al. (2021). We use the Document Sentiment Classification (DSC) dataset in the DIL setting. It consists of 10 subsets of product reviews with a positive or negative sentiment label. Each subset constitutes a separate task with $4,000$ training, $500$ validation, and $500$ test samples. Since the tasks are similar and only differ in the product domain, this dataset is used to evaluate knowledge transfer. In the CIL setup, we use the earlier described 20ng dataset (Lang, 1995). Similarly to Ke et al. (2021), we create a sequence of 10 tasks consisting of two classes each. This setup is mainly used to test the models' abilities to overcome catastrophic forgetting. For the TIL setup, we create a sequence

of tasks by combining four tasks from the GLUE benchmark (Wang et al., 2018). This dataset, which we call 4GLUE, includes four different task-types: The RTE dataset is used for testing natural language inference, the MRPC is used for measuring semantic textual similarity, the SST-2 is a popular dataset for sentiment analysis, and the QQP dataset which is used for natural language inference and question answering.

For the single-head CIL experiments, we use two different versions of the R21578 news dataset (Lewis, 1997), R8 (includes 8 classes) and R52 (with 52 classes). Due to the R21578 dataset's highly skewed class frequency distribution, we simulate a low-resource training scenario and include only 100 samples from each class in both datasets. On R8, we divide the dataset into 8 increments, with one class for each increment. On R52, we create 26 increments, each with two random classes.

**Procedure**  We follow the CL evaluation procedure of De Lange et al. (2021). A model is trained sequentially on all tasks and is evaluated by averaging the test performance of each task recorded after the final training increment. This results in each task in the sequence being a binary classification problem. In the multi-head configurations (for CIL and TIL), we use a separate decoder for each task and provide the task ID during training and evaluation. To further investigate the continual learning capabilities of the DKVB, we implement the single-head CIL setup of Träuble et al. (2023). Compared to the multi-head CIL task, this setup is considered to be more challenging and lacks explicit task boundaries. For its evaluation, the models are tested on the whole test data after each increment, including previously unseen classes.

**Baselines**  We use the best-performing methods reported in Ke et al. (2021), selecting one from each CL approach (cf. Section 2.1). From regularization-based methods, we choose **EWC** (Serra et al., 2018), a common baseline with strong performance in many CL studies. **DER++** (Buzzega et al., 2020) belongs to the replay-based methods and uses distilled knowledge from past experiences to guide the incremental training process. **OWM** (Zeng et al., 2019) is an optimization-based approach that constrains the gradient updates to a direction orthogonal to the input space of previously trained tasks. Lastly, **CTR** (Ke et al., 2021) is an architecture-based approach that utilizes capsule networks to

Table 2: Average Accuracy and Macro F1 scores of the tasks in the three continual learning scenarios, using the average and standard deviation (in subscript) of 5 runs with randomized sequence orders.

| CL Method | Model | DIL (DSC) | | CIL (20NG) | | TIL (4GLUE) | |
|---|---|---|---|---|---|---|---|
| | | Acc | F1 | Acc | F1 | Acc | F1 |
| NCL | BERT | $88.50_{0.63}$ | $87.83_{0.64}$ | $53.95_{0.49}$ | $38.94_{0.53}$ | $62.12_{0.58}$ | $58.29_{0.60}$ |
| NCL | BERT (frozen) | $87.42_{2.23}$ | $86.58_{2.17}$ | $96.35_{1.03}$ | $96.31_{1.22}$ | $71.90_{0.11}$ | $\mathbf{69.10}_{0.12}$ |
| NCL | Adapter-BERT | $88.71_{2.20}$ | $\mathbf{88.10}_{2.24}$ | $65.61_{9.34}$ | $58.62_{11.63}$ | $68.74_{0.30}$ | $63.86_{0.53}$ |
| DER++ | BERT (frozen) | $84.30_{1.54}$ | $82.85_{1.94}$ | $59.68_{9.23}$ | $47.62_{13.74}$ | $70.64_{3.27}$ | $68.68_{3.89}$ |
| EWC | BERT (frozen) | $86.21_{4.89}$ | $85.53_{4.98}$ | $\mathbf{96.80}_{0.20}$ | $\mathbf{96.80}_{0.20}$ | $66.54_{9.74}$ | $58.26_{14.84}$ |
| OWM | BERT (frozen) | $86.06_{2.63}$ | $85.28_{2.66}$ | $88.80_{0.28}$ | $88.16_{0.30}$ | $67.54_{2.11}$ | $61.90_{2.57}$ |
| CTR | Adapter-BERT | $\mathbf{88.73}_{0.35}$ | $87.98_{0.37}$ | $95.53_{0.14}$ | $95.52_{0.16}$ | $\mathbf{72.71}_{0.19}$ | $66.42_{0.78}$ |
| DKVB-NP Incremental | BERT (frozen) | $80.99_{2.07}$ | $79.58_{1.91}$ | $59.67_{1.59}$ | $54.58_{1.66}$ | $58.12_{2.89}$ | $50.14_{1.90}$ |
| DKVB-NP Oracle | BERT (frozen) | $\mathbf{83.93}_{1.11}$ | $\mathbf{81.98}_{1.74}$ | $\mathbf{97.06}_{0.22}$ | $95.84_{0.95}$ | $\mathbf{69.65}_{0.34}$ | $\mathbf{68.92}_{0.38}$ |
| DKVB-NP Generic | BERT (frozen) | $82.12_{0.20}$ | $80.97_{0.09}$ | $96.30_{0.07}$ | $\mathbf{96.27}_{0.10}$ | $68.79_{0.51}$ | $65.37_{0.03}$ |
| DKVB-P Incremental | BERT (frozen) | $74.09_{4.88}$ | $68.01_{5.10}$ | $57.81_{2.00}$ | $52.89_{2.77}$ | $58.77_{3.23}$ | $51.02_{1.81}$ |
| DKVB-P Oracle | BERT (frozen) | $81.18_{0.61}$ | $80.47_{0.52}$ | $95.22_{0.44}$ | $95.09_{0.25}$ | $58.65_{1.43}$ | $51.81_{1.53}$ |
| DKVB-P Generic | BERT (frozen) | $71.71_{1.30}$ | $57.75_{0.95}$ | $92.76_{0.88}$ | $92.73_{0.89}$ | $61.40_{0.57}$ | $54.76_{0.42}$ |

prevent catastrophic forgetting and facilitate knowledge transfer. We also include three baselines without any additional forgetting or knowledge transfer handling, noted as *naive continual learning* (**NCL**).

For DKVB we take the best-performing architectures from Section 3.3, and include both the parametric (**DKVB-P**) and non-parametric (**DKVB-NP**) variants. We experiment with three different strategies for key initialization. In the first two strategies, we use the training data for initializing the keys: in the incremental setup, the keys are optimized in a continual fashion before each task using only the training data of the given increment (denoted as **Incremental**), while in the full initialization setup, the keys are initialized once before training, using the full training input distribution (denoted as **Oracle**). In the third setup, we use a cross-domain corpus different from the training data to create general-purpose keys (denoted as **Generic**). For this, we use a small version of the English Wikipedia dump[2], which is commonly included in pre-training datasets. In all three setups, we use an EMA decay of $0.2$. For the Incremental and Oracle setups, the key initialization is set to three epochs, while for the Generic we use one epoch.

All CL methods (except CTR) are applied to a frozen BERT model and have a single-head configuration without any task-ID information for the DIL scenario and a multi-head configuration with task-ID provision on the CIL and TIL scenarios. CTR is based on an Adapter-BERT (Houlsby et al., 2019) backbone and requires a multi-head setup

and task-ID information for its dynamic architecture in all scenarios. For the single-head CIL experiments, we include the naive baselines and the replay-based DER++ method. The rest of the CL baseline methods either require explicit task boundaries for optimal performance (OWM, EWC) or only work in a multi-head configuration (CTR).

# 6 Results

**Main Experiments** The results of the main experiments are shown in Table 2. In the DIL setting the difference in accuracy between the baseline methods is low, with CTR having the highest score of $88.73\%$. The performance of the DKVB variants in this scenario is below the baselines. In the CIL setting, there is a substantial variation between model performance, with half of the CL methods achieving over 90% accuracy, while BERT NCL, DER++, and the incremental DKVB variants have an accuracy score below $60\%$. The best result on the CIL dataset was achieved with the non-parametric DKVB Oracle ($97.06\%$) followed by EWC ($96.80\%$) and BERT frozen NCL ($96.35\%$). In the TIL scenario, the highest accuracy scores were achieved with BERT frozen NCL ($71.90\%$) and CTR ($72.71\%$). Within the DKVB variants the best performance was consistently seen with the non-parametric Oracle variant, closely followed by the non-parametric Generic variant. On the CIL and TIL scenarios both of these methods outperformed most of the baselines. Additional measures of the backward transfer performances can be found in Appendix Section C.

---

[2]https://huggingface.co/datasets/wikipedia

**Runtime** We measure the average epoch runtimes for each model to compare the computational costs of the different methods. The results can be found in Table 3. Among the evaluated methods, DKVB achieves runtime closest to NCL with a frozen BERT, where training is limited to optimizing a parametric decoder. While the regularization-based EWC and the optimization-based OWM methods also achieve a runtime comparable to the NCL frozen BERT model, adding replay in DER++ and dynamic architecture in CTR substantially increases runtime. The key initialization process scales with the number of samples, but the overall computational cost of DKVB remains lower than most continual learning methods since initialization is done once before training and involves just a forward pass. The average runtime of key initialization is shown in Table 4

Table 3: Per-epoch training runtimes (in seconds), averaged over a single run. Standard deviations are shown as subscripts.

| CL Method | Model | DIL (DSC) | CIL (20NG) | TIL (4GLUE) |
|---|---|---|---|---|
| NCL | BERT | $20.6_{3.0}$ | $8.9_{0.0}$ | $482.3_{659.4}$ |
| NCL | BERT (frozen) | $4.4_{0.6}$ | $1.9_{0.0}$ | $105.5_{144.2}$ |
| NCL | Adapter-BERT | $24.1_{3.4}$ | $10.4_{0.0}$ | $566.2_{772.7}$ |
| DER++ | BERT (frozen) | $26.2_{0.9}$ | $7.4_{2.5}$ | $249.7_{361.2}$ |
| EWC | BERT (frozen) | $8.0_{0.08}$ | $2.3_{0.2}$ | $129.0_{176.6}$ |
| OWM | BERT (frozen) | $6.7_{0.3}$ | $2.0_{0.1}$ | $108.9_{148.6}$ |
| CTR | Adapter-BERT | $487.1_{0.4}$ | $195.2_{0.1}$ | $3011.7_{0.0}$ |
| DKVB-NP | BERT (frozen) | $4.67_{0.6}$ | $2.00_{0.0}$ | $109.35_{149.2}$ |
| DKVB-P | BERT (frozen) | $4.88_{0.7}$ | $2.07_{0.0}$ | $114.38_{156.5}$ |

Table 4: Per-epoch key initialization runtimes (in seconds) and corresponding sample sizes. Standard deviations are shown as subscripts.

| Key Initialization | DIL (DSC) | CIL (20NG) | TIL (4GLUE) |
|---|---|---|---|
| Incremental | $4.7_{0.6}$ (n=4 000) | $1.9_{0.0}$ (n=1 600) | $111.6_{152.5}$ (n=87 470) |
| Oracle | $46.9_{0.5}$ (n=40 000) | $19.5_{0.1}$ (n=16 000) | $535.22_{1.9}$ (n=349 881) |
| Generic | $469.0_{2.1}$ (n=205 328) | $469.0_{2.1}$ (n=205 328) | $469.0_{2.1}$ (n=205 328) |

**Single-head Class Incremental Learning** The single-head class incremental learning results are shown in Figure 2. The highest accuracy scores are 81.17% on R8 and 47.78% on R52. Both scores were achieved with the non-parametric DKVB variant using the Generic and Oracle key initialization, respectively. On both datasets, the non-DKVB models, which included the BERT frozen NCL and DER++, displayed sharp drops in performance between increments, indicating the occurrence of

catastrophic forgetting and overfitting on the current training increment. DER++ showcased better performance than the naive baseline but still underperformed the Oracle and Generic variants, with a final accuracy score of 16.70% on R8 and 35.75% on R52. The detailed results with additional models (BERT NCL, EWC) can be found in Appendix Section C.



Figure 2: Progressive test accuracy scores in the single-head class increment learning setup, averaged over 5 runs with fixed sequence order

## 7 Discussion

**Key Insights** Our experiments show that fine-tuning encoder-only language models with an optimal discrete key-value bottleneck architecture achieves comparable results to partial fine-tuning in standard learning scenarios, but greatly benefits CL, both in terms of performance and efficiency. The best key initialization is obtained by unsupervised access to the full input feature distribution, but utilizing a general-purpose corpus for key initialization is also a viable option for NLP tasks. Below, we discuss these key insights.

**Architectural Variants** We found that employing pooling before the bottleneck has a substantial negative effect on the model performance (see Section 3.3). This suggests that in contrast to lower dimensional vision tasks (Träuble et al., 2023), it is necessary to retain the full dimensionality of the

text encodings. Similarly, CLS pooling is inferior compared to mean pooling across all setups. Segmenting on the token dimension only worked in the case of parametric decoding, indicating that the DKVB module's output yields better representational power if the segmentation happens on the hidden dimension. An additional linear layer after the DKVB module, acting as a parametric decoder, can compensate for encodings with weaker representational power. However, a non-parametric decoder produces comparable results in most configurations.

**Continual Learning**  In the CL experiments reported in Section 5, the non-parametric DKVB variants achieved comparable results to other CL methods and maintained runtimes on par with the NCL frozen BERT variant, but only when using pre-initialized keys.

When using incremental key initialization, performance was consistently subpar, indicating that DKVB requires access to a general-purpose corpus or the full input distribution to achieve competitive performance. While having access to the full data distribution may be unrealistic in practice, our experiments show that this is not needed in NLP. Rather, when initializing the keys using a general-purpose corpus, we obtain results that are close to the Oracle setup.

The largest performance drop between the DKVB variants and other CL methods was seen in the DIL setting. This suggests that DKVB's strength in preventing catastrophic forgetting through distinct key-value bindings becomes its weakness in DIL, as this compartmentalization restricts the model's ability to transfer knowledge across different domains. Notably, NCL methods achieve similar results as the CL methods in this DIL setup, indicating that pre-trained language models without a bottleneck are already well-suited for domain incremental learning.

In the CIL and TIL tasks, only the frozen BERT NCL variant showcased performance comparable to that of the CL methods. The strong performance of frozen BERT in these experiments suggests that if task-ID is available during testing, a simple multi-head configuration with a frozen encoder is often sufficient. Experiments in the single-head CIL setup have shown to be more challenging. As the models are tested on the full test set after each increment, ideally, they should exhibit a progressive increase in accuracy. But when no task-ID is pro-

vided and decoding is done with a single head, most models overfit and suffer catastrophic forgetting between increments, with DKVB being the only model to demonstrate improved CL capability in this scenario. This suggests that DKVB's unique architecture effectively maintains knowledge across tasks without needing task-specific heads.

## 8   Conclusion

The discrete key-value bottleneck offers an efficient approach to continual learning. It enables context-dependent updates in the model without explicit parameter isolation or dynamically expanding the architecture. Considering the special challenges of continual learning with text embeddings, we analyzed twelve architectural variants of the bottleneck. The best variants apply mean pooling after the bottleneck and utilize the hidden dimension of the encoded input representation to create the bottleneck heads.

We conducted a comprehensive evaluation across different continual learning settings in NLP, i.e., domain-incremental learning, class-incremental learning, and task-type incremental learning, and showed that with proper key initialization, the discrete key-value bottleneck offers consistent improvement in most settings and is comparable to dedicated continual learning methods from the literature. Moreover, we showed that it can be used even in the most challenging single-head continual learning scenarios when no task-ID is provided.

## 9   Limitations

Our study focuses solely on encoder-only language models. While this raises questions about whether our results could generalize to other model architectures, our choice was motivated by their preference for supervised fine-tuning scenarios where the balance between performance and computational efficiency is crucial. Our experiments were also limited to fine-tuning for classification-based downstream tasks. Consequently, it remains to be investigated whether our results extend to other NLP tasks, such as semantic search, entity extraction, or machine translation.

# References

Magdalena Biesialska, Katarzyna Biesialska, and Marta R Costa-jussà. 2020. Continual lifelong learning in natural language processing: A survey. In *Proceedings of the 28th International Conference on Computational Linguistics*, pages 6523–6541.

Zalán Borsos, Raphaël Marinier, Damien Vincent, Eugene Kharitonov, Olivier Pietquin, Matt Sharifi, Dominik Roblek, Olivier Teboul, David Grangier, Marco Tagliasacchi, et al. 2023. AudioLM: a language modeling approach to audio generation. *IEEE/ACM Transactions on Audio, Speech, and Language Processing*, 31:2523–2533.

Tom Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, et al. 2020. Language models are few-shot learners. *Advances in neural information processing systems*, 33:1877–1901.

Pietro Buzzega, Matteo Boschini, Angelo Porrello, Davide Abati, and Simone Calderara. 2020. Dark experience for general continual learning: a strong, simple baseline. *Advances in neural information processing systems*.

Arslan Chaudhry, Marc'Aurelio Ranzato, Marcus Rohrbach, and Mohamed Elhoseiny. 2018. Efficient lifelong learning with a-gem. In *International Conference on Learning Representations*.

Michael Crawshaw. 2020. Multi-task learning with deep neural networks: A survey. *arXiv preprint arXiv:2009.09796*.

Matthias De Lange, Rahaf Aljundi, Marc Masana, Sarah Parisot, Xu Jia, Aleš Leonardis, Gregory Slabaugh, and Tinne Tuytelaars. 2021. A continual learning survey: Defying forgetting in classification tasks. *IEEE transactions on pattern analysis and machine intelligence*, 44(7):3366–3385.

Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. 2018. BERT: Pre-training of deep bidirectional transformers for language understanding. *arXiv preprint arXiv:1810.04805*.

Lukas Galke and Ansgar Scherp. 2022. Bag-of-words vs. graph vs. sequence in text classification: Questioning the necessity of text-graphs and the surprising strength of a wide mlp. In *Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 4038–4051.

Anirudh Goyal, Alessandro Sordoni, Marc-Alexandre Côté, Nan Rosemary Ke, and Yoshua Bengio. 2017. Z-forcing: Training stochastic recurrent networks. *Advances in neural information processing systems*, 30.

Yiduo Guo, Wenpeng Hu, Dongyan Zhao, and Bing Liu. 2022. Adaptive orthogonal projection for batch and online continual learning. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 36, pages 6783–6791.

Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. 2015. Spatial pyramid pooling in deep convolutional networks for visual recognition. *IEEE transactions on pattern analysis and machine intelligence*, 37(9):1904–1916.

Neil Houlsby, Andrei Giurgiu, Stanislaw Jastrzebski, Bruna Morrone, Quentin De Laroussilhe, Andrea Gesmundo, Mona Attariyan, and Sylvain Gelly. 2019. Parameter-efficient transfer learning for NLP. In *International conference on machine learning*, volume 97, pages 2790–2799.

Ching-Yi Hung, Cheng-Hao Tu, Cheng-En Wu, Chien-Hung Chen, Yi-Ming Chan, and Chu-Song Chen. 2019. Compacting, picking and growing for unforgetting continual learning. *Advances in Neural Information Processing Systems*, 32.

Dieuwke Hupkes, Mario Giulianelli, Verna Dankers, Mikel Artetxe, Yanai Elazar, Tiago Pimentel, Christos Christodoulopoulos, Karim Lasri, Naomi Saphra, Arabella Sinclair, et al. 2023. A taxonomy and review of generalization research in NLP. *Nature Machine Intelligence*.

Khurram Javed and Martha White. 2019. Meta-learning representations for continual learning. *Advances in neural information processing systems*, 32.

Zixuan Ke and Bing Liu. 2022. Continual learning of natural language processing tasks: A survey. *arXiv preprint arXiv:2211.12701*.

Zixuan Ke, Bing Liu, Nianzu Ma, Hu Xu, and Lei Shu. 2021. Achieving forgetting prevention and knowledge transfer in continual learning. *Advances in Neural Information Processing Systems*, 34:22443–22456.

Gyuhak Kim, Changnan Xiao, Tatsuya Konishi, Zixuan Ke, and Bing Liu. 2022. A theoretical study on solving continual learning. *Advances in neural information processing systems*, 35:5065–5079.

James Kirkpatrick, Razvan Pascanu, Neil Rabinowitz, Joel Veness, Guillaume Desjardins, Andrei A Rusu, Kieran Milan, John Quan, Tiago Ramalho, Agnieszka Grabska-Barwinska, et al. 2017. Overcoming catastrophic forgetting in neural networks. *Proceedings of the national academy of sciences*, 114(13):3521–3526.

Ken Lang. 1995. Newsweeder: Learning to filter netnews. In *International Conference on Machine Learning*, pages 331–339.

David Lewis. 1997. Reuters-21578 Text Categorization Collection. UCI Machine Learning Repository. DOI: https://doi.org/10.24432/C52G6M.

Haitao Li, Qingyao Ai, Jia Chen, Qian Dong, Zhijing Wu, and Yiqun Liu. 2025. Blade: Enhancing black-box large language models with small domain-specific models. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 39, pages 24422–24430.

Dianbo Liu, Alex Lamb, Xu Ji, Pascal Junior Tikeng Notsawo, Michael Mozer, Yoshua Bengio, and Kenji Kawaguchi. 2023. Adaptive discrete communication bottlenecks with dynamic vector quantization for heterogeneous representational coarseness. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 37, pages 8825–8833.

Dianbo Liu, Alex M Lamb, Kenji Kawaguchi, Goyal Anirudh, Chen Sun, Michael C Mozer, and Yoshua Bengio. 2021. Discrete-valued neural communication. *Advances in Neural Information Processing Systems*, 34:2109–2121.

Yinhan Liu, Myle Ott, Naman Goyal, Jingfei Du, Mandar Joshi, Danqi Chen, Omer Levy, Mike Lewis, Luke Zettlemoyer, and Veselin Stoyanov. 2019. RoBERTa: A robustly optimized BERT pretraining approach. *arXiv preprint arXiv:1907.11692*.

David Lopez-Paz and Marc'Aurelio Ranzato. 2017. Gradient episodic memory for continual learning. *Advances in neural information processing systems*, 30.

Yun Luo, Zhen Yang, Xuefeng Bai, Fandong Meng, Jie Zhou, and Yue Zhang. 2023. Investigating forgetting in pre-trained representations through continual learning. *arXiv preprint arXiv:2305.05968*.

Sergio Naval Marimont and Giacomo Tarroni. 2021. Anomaly detection through latent space restoration using vector quantized variational autoencoders. In *2021 IEEE 18th International Symposium on Biomedical Imaging (ISBI)*, pages 1764–1767. IEEE.

Michael McCloskey and Neal J Cohen. 1989. Catastrophic interference in connectionist networks: The sequential learning problem. In *Psychology of learning and motivation*, volume 24, pages 109–165. Elsevier.

Seyed Iman Mirzadeh, Mehrdad Farajtabar, Dilan Gorur, Razvan Pascanu, and Hassan Ghasemzadeh. 2020. Linear mode connectivity in multitask and continual learning. *arXiv preprint arXiv:2010.04495*.

Saleh Momeni, Sahisnu Mazumder, and Bing Liu. 2025. Continual learning using a kernel-based method over foundation models. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 39, pages 19528–19536.

Niklas Muennighoff, Nouamane Tazi, Loic Magne, and Nils Reimers. 2023. Mteb: Massive text embedding benchmark. In *Proceedings of the 17th Conference of the European Chapter of the Association for Computational Linguistics*, pages 2014–2037.

Haechan Noh, Sangeek Hyun, Woojin Jeong, Hanshin Lim, and Jae-Pil Heo. 2023. Disentangled representation learning for unsupervised neural quantization. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 12001–12010.

Oleksiy Ostapenko, Timothee Lesort, Pau Rodriguez, Md Rifat Arefin, Arthur Douillard, Irina Rish, and Laurent Charlin. 2022. Continual learning with foundation models: An empirical study of latent replay. In *Conference on Lifelong Learning Agents*, pages 60–91. PMLR.

Muhammad Reza Qorib, Geonsik Moon, and Hwee Tou Ng. 2024. Are decoder-only language models better than encoder-only language models in understanding word meaning? In *Annual Meeting of the Association for Computational Linguistics*.

Anastasia Razdaibiedina, Yuning Mao, Rui Hou, Madian Khabsa, Mike Lewis, and Amjad Almahairi. 2023. Progressive prompts: Continual learning for language models. In *The Eleventh International Conference on Learning Representations*.

Victor Sanh, Lysandre Debut, Julien Chaumond, and Thomas Wolf. 2019. DistilBERT, a distilled version of BERT: smaller, faster, cheaper and lighter. *arXiv preprint arXiv:1910.01108*.

Thomas Scialom, Tuhin Chakrabarty, and Smaranda Muresan. 2022. Fine-tuned language models are continual learners. In *Proceedings of the 2022 Conference on Empirical Methods in Natural Language Processing*, pages 6107–6122.

Joan Serra, Didac Suris, Marius Miron, and Alexandros Karatzoglou. 2018. Overcoming catastrophic forgetting with hard attention to the task. In *International conference on machine learning*, pages 4548–4557. PMLR.

Vedant Shah, Frederik Träuble, Ashish Malik, Hugo Larochelle, Michael Mozer, Sanjeev Arora, Yoshua Bengio, and Anirudh Goyal. 2023. Unlearning via sparse representations. *arXiv preprint arXiv:2311.15268*.

Haizhou Shi, Zihao Xu, Hengyi Wang, Weiyi Qin, Wenyuan Wang, Yibin Wang, Zifeng Wang, Sayna Ebrahimi, and Hao Wang. 2024. Continual learning of large language models: A comprehensive survey. *ACM Computing Surveys*.

Hanul Shin, Jung Kwon Lee, Jaehong Kim, and Jiwon Kim. 2017. Continual learning with deep generative replay. *Advances in neural information processing systems*, 30.

Frederik Träuble, Anirudh Goyal, Nasim Rahaman, Michael Curtis Mozer, Kenji Kawaguchi, Yoshua Bengio, and Bernhard Schölkopf. 2023. Discrete key-value bottleneck. In *International Conference on Machine Learning*, pages 34431–34455. PMLR.

Gido M Van de Ven and Andreas S Tolias. 2019. Three scenarios for continual learning. *arXiv preprint arXiv:1904.07734*.

Aaron Van Den Oord, Oriol Vinyals, et al. 2017. Neural discrete representation learning. *Advances in neural information processing systems*, 30.

Alex Wang, Amanpreet Singh, Julian Michael, Felix Hill, Omer Levy, and Samuel R Bowman. 2018. Glue: A multi-task benchmark and analysis platform for natural language understanding. *EMNLP 2018*, page 353.

Liyuan Wang, Xingxing Zhang, Hang Su, and Jun Zhu. 2024. A comprehensive survey of continual learning: theory, method and application. *IEEE Transactions on Pattern Analysis and Machine Intelligence*.

Mitchell Wortsman, Vivek Ramanujan, Rosanne Liu, Aniruddha Kembhavi, Mohammad Rastegari, Jason Yosinski, and Ali Farhadi. 2020. Supermasks in superposition. *Advances in Neural Information Processing Systems*, 33:15173–15184.

Wilson Yan, Yunzhi Zhang, Pieter Abbeel, and Aravind Srinivas. 2021. VideoGPT: Video generation using VQ-VAE and transformers. *arXiv preprint arXiv:2104.10157*.

Wenpeng Yin, Jia Li, and Caiming Xiong. 2022. Contintin: Continual learning from task instructions. In *Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 3062–3072.

Hao Yu, Zachary Yang, Kellin Pelrine, Jean Francois Godbout, and Reihaneh Rabbany. 2023. Open, closed, or small language models for text classification? *arXiv preprint arXiv:2308.10092*.

Lifan Yuan, Yangyi Chen, Ganqu Cui, Hongcheng Gao, Fangyuan Zou, Xingyi Cheng, Heng Ji, Zhiyuan Liu, and Maosong Sun. 2023. Revisiting out-of-distribution robustness in nlp: Benchmarks, analysis, and llms evaluations. *Advances in Neural Information Processing Systems*, 36:58478–58507.

Guanxiong Zeng, Yang Chen, Bo Cui, and Shan Yu. 2019. Continual learning of context-dependent processing in neural networks. *Nature Machine Intelligence*, 1(8):364–372.

## Supplementary Materials

## A  Extended Related Work

**Regularization-based methods**  This family of methods involves incorporating explicit regularization terms to maintain a balance between the old and new tasks. This is usually done by adding penalty or regularization to the loss function to prevent large changes to parameters deemed important for old tasks (Wang et al., 2024). A popular method in this family is Elastic Weight Consideration (EWC), which calculates the importance of parameters with the Fisher information matrix, and applies smaller updates to weights deemed critical for earlier tasks (Kirkpatrick et al., 2017).

**Replay-based methods**  These methods either retain a subset of training examples from previous tasks in memory such as A-GEM (Chaudhry et al., 2018), or learn to generate pseudo samples from previous tasks, like in DGR (Shin et al., 2017). These samples are then incorporated into the training regimen of new tasks. While this can alleviate catasthropical forgetting, the size of a memory buffer is limited, which can potentially affect generalizability (Wang et al., 2024).

**Optimization-based methods**  Explicitly manipulating the optimization process is another way to tackle the challenges of continual learning. Gradient-projection methods ensure that gradient updates happen exclusively in the orthogonal direction to the gradients of an old tasks, thereby preventing any impact on weights important for old tasks (Zeng et al., 2019; Guo et al., 2022). Meta learning strategies and methods focusing on obtaining flat minima in the loss landscape can be also utilized in continual learning (Javed and White, 2019; Mirzadeh et al., 2020).

**Architecture-based methods**  Methods in this family can be generally divided into *parameter isolation* and, *dynamic architecture* approaches, depending on whether the model architecture is fixed or not (Wang et al., 2024). Models such as Sup-Sup (Wortsman et al., 2020) and HAT (Serra et al., 2018) optimize a binary mask to selectively choose dedicated parameters for each task and fall under the parameter isolation category. Other methods dynamically expand the model with new parameters to increase capacity for learning new tasks (Ke et al., 2021; Hung et al., 2019).

**Instruction-based methods**  This family is unique to the field of NLP. These methods are based on task-specific instructions given to encoder-decoder or decoder only language models when a new task is encountered. While some methods in this family show promising knowledge transfer capabilities (Scialom et al., 2022; Yin et al., 2022; Razdaibiedina et al., 2023), without explicit fine-tuning they are mostly limited by the knowledge acquired in the pre-training phase.

## B  Extended Experimental Setup

**Implementation**  For all model backbones in our experiments, we use the BERT-base model from Huggingface[3] and use cross-entropy loss as our objective function. We base our discrete-key-value bottleneck implementation on the *vector-quantize-pytorch*[4] package. In the pre-experiments, we truncate each input sample to 256 tokens. For the main continual learning experiments, we rely on the *Py-Continual*[5] framework and reuse their implementations and hyperparameters on the baseline methods. To remain comparable to other studies using the *PyContinual* framework, we kept the default pre-processing steps, used a maximum token length of 128, and applied the default convolutional decoder of the baseline models. For the single-head class incremental learning experiments we use a fixed randomized sequence order when creating the increments, and used a token length of 256. The source code for our experiments alongside the models can be found at github.com/drndr/dkvb_nlp.

**Optimization**  As part of our pre-experiments, we also conducted a hyperparameter search and a sensitivity analysis on the bottleneck parameters. Outside of the selected hyperparameters and bottleneck parameters, all other configurations remained fixed during the search. Our experiments use the BERT-base architecture with a hidden size of 768. For the optimizer, we chose AdamW with a weight decay of 0.01. The dropout rate for the parametric decoder was set to 0.1. For the reference fully fine-tuned BERT numbers we reused the hyperparameters reported in (Galke and Scherp, 2022), for the frozen BERT variant we relied on the parametric DKVB variant hyperparameters with mean pooling. During fine-tuning, we carefully

---

[3]https://huggingface.co/bert-base-uncased
[4]https://github.com/lucidrains/vector-quantize-pytorch
[5]https://github.com/ZixuanKe/PyContinual

optimized the models on both datasets using grid-search-based manual tuning. A search space for the selected hyperparameters was defined, specifically we chose the batch size from $\{8, 16, 32\}$, the number of epochs from $\{5, 10\}$, the learning rate for the values layer from $\{1e\text{-}1, 1e\text{-}2, 1e\text{-}3\}$, and the decoder learning rate from $\{1e\text{-}3, 1e\text{-}4, 1e\text{-}5\}$. The best performing (based on the validation loss) configurations for each architecture variant can be seen in Table 5. The hyperparameters were reused for the continual learning main experiments.

In the single-head class incremental learning experiments we conducted an additional manual hyperparameter search for the DKVB variants and found a batch size of 16 with a global learning rate of $1e-2$ to be the best performing. For EWC we set the lambda parameter to 5,000. In the DER++ model we used a memory buffer size of 256 and set the sampling reate in each increment to 16.

## B.1 Bottleneck Parameters

In our experiments, we rely on the optimal bottleneck parameter analysis of (Träuble et al., 2023). Additionally, we also conduct a small sensitivity study for the discrete key dimension and number of key-value pairs on the R8 dataset. For this, we use the DKVB-NP model variant from the pre-experiments and keep everything fixed, changing only these two bottleneck parameters. For the base hyperparameters, we reuse the best-performing configurations.

**Key dimension** The number of dimensions of the discrete keys strongly influences the utility of the bottleneck. This can be explained as follows. Keys that have too few dimensions increase the chance of unintended key sharing between inputs from different distributions, while discrete keys with too high dimensionality can lead to insufficient coverage of the embedding space. Similarly to (Träuble et al., 2023) we found the optimal key dimension to be between 8 to 12. The results of this analysis are depicted in Figure 3.

**Number of key-value pairs** The number of key-value pairs determines the size of the discretized representational space. In accordance with the analysis of (Träuble et al., 2023), we found that increasing the number of key-value pairs leads to a performance increase. Eventually further increments no longer yield substantial improvements in performance. Note that increasing this parameter also leads to increased model size and increases the



(a)



(b)

Figure 3: Assessing the sensitivity of bottleneck parameters in regards of test accuracy: (a) Dimensionality of discrete key (b) Number of key-value pairs

computational costs of key initialization as well. The results of this analysis are depicted in Figure 3.

## C Extended Results

**Architectural Variants** Results obtained using the RoBERTa (Table 6) and DistilBERT (Table 7) models demonstrate comparable performance patterns to those observed with BERT. Interestingly, DistilBERT produced slightly better accuracy on most architecture variants compared to the other two models. However the highest performance was consistently seen with mean pooling after the bottleneck on all three models. DistilBERT's improved performance can likely be attributed to differences in its tokenization and pooling implementation compared to other models.

**Continual Learning Experiments** In the field of continual learning additional metrics are often used to measure the performance over incremental learning. Two often used metrics are Forward Transfer (FWT) and Backward Transfer (BWT) (Lopez-Paz and Ranzato, 2017). BWT refers to how learning a new task affects performance on a previously learned task. It can be positive, when learning a new task improves performance on the earlier task, or negative, when it worsens it. Severe negative backward transfer is often called catas-

Table 5: Best hyperparameter configuration for each architecture variant based on validation loss

| Model | Segmentation | Pooling | (A) Dataset: R8 | | | | (B) Dataset: 20ng | | | |
|---|---|---|---|---|---|---|---|---|---|---|
| | | | #Epoch | Batch size | Values LR | Decoder LR | #Epoch | Batch size | Values LR | Decoder LR |
| DKVB-P | hidden | Before CLS | 5 | 32 | 1e-2 | 1e-4 | 10 | 16 | 1e-1 | 1e-3 |
| DKVB-P | hidden | Before Mean | 5 | 32 | 1e-1 | 1e-4 | 10 | 16 | 1e-2 | 1e-4 |
| DKVB-P | hidden | After CLS | 10 | 16 | 1e-2 | 1e-4 | 5 | 16 | 1e-2 | 1e-4 |
| DKVB-P | hidden | After Mean | 10 | 16 | 1e-2 | 1e-3 | 5 | 16 | 1e-2 | 1e-3 |
| DKVB-P | token | After CLS | 10 | 16 | 1e-2 | 1e-3 | 10 | 16 | 1e-2 | 1e-3 |
| DKVB-P | token | After Mean | 10 | 16 | 1e-2 | 1e-3 | 10 | 16 | 1e-2 | 1e-3 |
| DKVB-NP | hidden | Before CLS | 5 | 32 | 1e-1 | - | 5 | 32 | 1e-1 | - |
| DKVB-NP | hidden | Before Mean | 5 | 32 | 1e-1 | - | 10 | 32 | 1e-1 | - |
| DKVB-NP | hidden | After CLS | 10 | 16 | 1e-1 | - | 5 | 16 | 1e-1 | - |
| DKVB-NP | hidden | After Mean | 10 | 32 | 1e-1 | - | 10 | 16 | 1e-1 | - |
| DKVB-NP | token | After CLS | 10 | 32 | 1e-1 | - | 10 | 32 | 1e-2 | - |
| DKVB-NP | token | After Mean | 10 | 32 | 1e-1 | - | 10 | 32 | 1e-2 | - |

Table 6: Accuracy and standard deviation (in subscript) of the different DKVB architecture variants with RoBERTa on the R8 and 20ng datasets in a non-continual, standard learning setup, averaged over 5 runs.

| Decoder | Segmentation | Dataset: R8 | | | | Dataset: 20ng | | | |
|---|---|---|---|---|---|---|---|---|---|
| | | Pooling Before | | Pooling After | | Pooling Before | | Pooling After | |
| | | CLS | Mean | CLS | Mean | CLS | Mean | CLS | Mean |
| Parametric | hidden | $49.48_{0.05}$ | $91.36_{0.40}$ | $91.73_{0.81}$ | $94.25_{0.26}$ | $51.05_{0.43}$ | $56.63_{0.39}$ | $52.18_{1.11}$ | $\mathbf{75.08}_{0.21}$ |
| | token | - | - | $90.02_{1.05}$ | $94.05_{0.31}$ | - | - | $19.86_{1.03}$ | $27.30_{0.92}$ |
| Non Parametric | hidden | $49.45_{0.02}$ | $92.05_{0.29}$ | $74.53_{1.74}$ | $93.04_{0.20}$ | $56.83_{0.35}$ | $60.42_{0.35}$ | $53.10_{1.37}$ | $70.33_{0.78}$ |
| | token | - | - | $58.74_{0.92}$ | $66.33_{0.74}$ | - | - | $9.89_{0.66}$ | $12.51_{1.67}$ |
| RoBERTa (frozen) w/o DKVB | | $94.29_{0.17}$ | | | | $69.42_{0.30}$ | | | |
| RoBERTa w/o DKVB | | $97.54_{0.51}$ | | | | $83.36_{0.30}$ | | | |

Table 7: Accuracy and standard deviation (in subscript) of the different DKVB architecture variants with DistilBERT on the R8 and 20ng datasets in a non-continual, standard learning setup, averaged over 5 runs.

| Decoder | Segmentation | Dataset: R8 | | | | Dataset: 20ng | | | |
|---|---|---|---|---|---|---|---|---|---|
| | | Pooling Before | | Pooling After | | Pooling Before | | Pooling After | |
| | | CLS | Mean | CLS | Mean | CLS | Mean | CLS | Mean |
| Parametric | hidden | $90.22_{0.30}$ | $92.17_{0.25}$ | $90.46_{0.44}$ | $\mathbf{96.09}_{0.21}$ | $56.26_{0.53}$ | $60.24_{0.20}$ | $60.26_{0.31}$ | $\mathbf{79.79}_{0.49}$ |
| | token | - | - | $89.24_{0.93}$ | $94.78_{0.86}$ | - | - | $45.45_{1.13}$ | $68.06_{0.90}$ |
| Non Parametric | hidden | $89.79_{0.24}$ | $92.02_{0.38}$ | $90.92_{0.22}$ | $95.09_{0.27}$ | $60.73_{0.48}$ | $59.98_{0.30}$ | $61.43_{0.50}$ | $75.11_{0.44}$ |
| | token | - | - | $66.37_{0.49}$ | $72.65_{0.24}$ | - | - | $12.04_{0.51}$ | $18.70_{1.03}$ |
| DistilBERT (frozen) w/o DKVB | | $94.62_{0.16}$ | | | | $68.56_{0.38}$ | | | |
| DistilBERT w/o DKVB | | $97.83_{0.24}$ | | | | $83.84_{0.23}$ | | | |

trophic forgetting. FWT describes how learning a new task influences performance on a future task. Since pre-trained language models already posses high transfer learning capabilities (Brown et al., 2020), its difficult to isolate the effect of learning specific task on future performance. Therefore we focus on BWT which is formally defined as:

$$\text{BWT} = \frac{1}{T-1} \sum_{i=1}^{T-1} (R_{T,i} - R_{i,i}) \qquad (1)$$

where $R \in \mathbb{R}^{T \times T}$ is the results matrix of an incremental learning scenario with $T$ tasks, where each entry $R_{i,j}$ being the test accuracy on task $j$ after training on task $i$ (Lopez-Paz and Ranzato, 2017). We report the BWT numbers on the three continual learning scenarios in Table 8.

Table 8: Average Backward Transfer (BWT) scores on the three continual learning scenarios

| CL Method | Model | DIL (DSC) | CIL (20NG) | TIL (4GLUE) |
|---|---|---|---|---|
| NCL | BERT | 0.29 | −29.77 | −20.00 |
| NCL | BERT (frozen) | −0.10 | −0.38 | −6.70 |
| NCL | Adapter-BERT | 0.39 | −20.01 | −16.05 |
| DER++ | BERT (frozen) | −0.58 | −27.11 | −7.05 |
| EWC | BERT (frozen) | 0.06 | **-0.27** | −10.84 |
| OWM | BERT (frozen) | 0.26 | −15.44 | −8.23 |
| CTR | Adapter-BERT | **0.49** | −0.50 | **-6.19** |
| DKVB-NP Incremental | BERT (frozen) | −3.06 | −27.73 | −21.30 |
| DKVB-NP Oracle | BERT (frozen) | **-0.88** | **-0.12** | **-7.22** |
| DKVB-NP Generic | BERT (frozen) | −1.17 | −0.29 | −7.97 |
| DKVB-P Incremental | BERT (frozen) | −4.88 | −29.05 | −20.99 |
| DKVB-P Oracle | BERT (frozen) | −1.02 | −0.41 | −20.56 |
| DKVB-P Generic | BERT (frozen) | −6.24 | −4.94 | −16.00 |

Table 9: Mean accuracy scores of single-head class incremental learning experiments on R8, averaged over 5 runs with fixed sequence order

| Increment | # Test Samples | BERT | BERT-frozen | BERT-frozen DER++ | BERT-frozen EWC | DKVB-NP Incremental | DKVB-NP Oracle | DKVB-NP Wiki |
|---|---|---|---|---|---|---|---|---|
| 1. | 1596 | 31.79 | 31.79 | 31.79 | 31.79 | 31.79 | 31.79 | 31.79 |
| 2. | 253 | 5.52 | 5.52 | 8.60 | 31.79 | 5.57 | 31.80 | 31.79 |
| 3. | 2840 | 49.47 | 49.47 | 46.47 | 12.65 | 49.70 | 67.70 | 43.19 |
| 4. | 41 | 0.45 | 0.45 | 4.99 | 49.61 | 27.40 | 76.98 | 75.83 |
| 5. | 190 | 3.70 | 3.70 | 14.57 | 49.56 | 41.02 | 77.57 | 77.56 |
| 6. | 206 | 3.97 | 3.97 | 23.11 | 3.70 | 44.86 | 79.26 | 79.76 |
| 7. | 108 | 1.64 | 1.64 | 14.89 | 3,74 | 43.85 | 79.72 | 80.61 |
| 8. | 251 | 3.42 | 3.42 | 16.71 | 2.64 | 29.83 | 80.86 | 81.90 |

Table 10: Mean accuracy scores of single-head class incremental learning experiments on R52, averaged over 5 runs with fixed sequence order

| Increment | # Test Samples | BERT | BERT-frozen | BERT-frozen DER++ | BERT-frozen EWC | DKVB-NP Incremental | DKVB-NP Oracle | DKVB-NP Wiki |
|---|---|---|---|---|---|---|---|---|
| 1. | 45 | 0.70 | 0.50 | 0.56 | 0.50 | 0.73 | 0.54 | 0.60 |
| 2. | 1600 | 26.20 | 27.10 | 26.07 | 0.50 | 27.10 | 27.10 | 0.92 |
| 3. | 52 | 0.46 | 0.23 | 14.99 | 0.85 | 0.54 | 27.29 | 24.71 |
| 4. | 29 | 0.35 | 0.35 | 9.60 | 3.69 | 0.46 | 29.51 | 27.30 |
| 5. | 321 | 2.92 | 2.95 | 4.08 | 2.57 | 3.58 | 25.60 | 27.41 |
| 6. | 37 | 0.35 | 0.07 | 18.76 | 0.35 | 0.42 | 27.22 | 27.63 |
| 7. | 17 | 0.35 | 0.42 | 7.73 | 0.35 | 0.35 | 28.85 | 27.55 |
| 8. | 44 | 0.54 | 0.46 | 4.53 | 0.35 | 0.70 | 27.57 | 27.67 |
| 9. | 28 | 0.50 | 0.35 | 13.13 | 0.42 | 0.46 | 27.10 | 28.10 |
| 10. | 110 | 1.40 | 1.40 | 12.73 | 0.35 | 1.40 | 25.58 | 28.56 |
| 11. | 3046 | 42.17 | 42.17 | 44.18 | 0.35 | 44.82 | 42.17 | 28.87 |
| 12. | 16 | 0.97 | 0.35 | 19.78 | 0.70 | 0.42 | 43.71 | 29.17 |
| 13. | 10 | 0.15 | 0.23 | 25.31 | 0.97 | 0.35 | 44.85 | 29.48 |
| 14. | 193 | 3.15 | 3.15 | 21.58 | 0.42 | 3.15 | 52.64 | 30.30 |
| 15. | 213 | 3.38 | 3.38 | 23.62 | 3.38 | 2.95 | 42.52 | 31.80 |
| 16. | 154 | 1.09 | 1.09 | 6.43 | 3.38 | 1.55 | 42.83 | 35.83 |
| 17. | 145 | 1.40 | 1.40 | 26.11 | 1.47 | 1.83 | 45.40 | 36.12 |
| 18. | 32 | 0.50 | 0.50 | 23.19 | 0.50 | 0.50 | 45.52 | 17.95 |
| 19. | 203 | 3.15 | 3.15 | 31.43 | 1.83 | 3.30 | 45.71 | 11.10 |
| 20. | 227 | 3.15 | 3.15 | 25.42 | 3.15 | 3.62 | 38.94 | 10.34 |
| 21. | 2948 | 42.17 | 42.17 | 45.67 | 3.15 | 42.25 | 44.74 | 40.07 |
| 22. | 255 | 4.71 | 4.71 | 47.63 | 42.25 | 4.71 | 48.84 | 41.96 |
| 23. | 59 | 0.58 | 0.38 | 23.27 | 42.17 | 0.77 | 47.15 | 42.23 |
| 24. | 48 | 0.58 | 0.58 | 22.85 | 0.58 | 0.62 | 37.96 | 42.48 |
| 25. | 59 | 0.58 | 0.58 | 42.04 | 0.58 | 0.70 | 38.55 | 43.11 |
| 26. | 243 | 3.38 | 3.38 | 35.75 | 0.58 | 3.62 | 47.78 | 45.04 |

# Domain adaptation and question-answer pooling for Aphasia modelling

**Uwe Reichel**[1], **Monica Gonzalez Machorro**[1,2], **Lisa M. Ehlen**[3], **Pascal Hecker**[1,4],
**Dorothea Peitz**[3], **Cornelius Werner**[3,5], **Felix Burkhardt**[1,6], **Christian Kohlschein**[7],
**Florian Eyben**[1,8], **Björn Schuller**[1,2,9,10]

[1]audEERING GmbH, [2]Chair of Health Informatics, TUM University Hospital,
[3]Dep. of Neurology, Medical Faculty, RWTH Aachen University, [4]Hasso-Plattner Institute,
[5]Dep. of Neurology and Geriatrics, Johanniter Hospital Stendal, [6]TU Berlin,
[7]Accenture GmbH, [8]Agile Robots, [9]Imperial College, [10]Munich Center for Machine Learning
ureichel@audeering.com

## Abstract

In this study, we examine the impact of domain adaptation and question-answer pooling on text-based aphasia prediction with standard and medically specialised BERT models for a German corpus. Modelling tasks comprise aphasia type classification as well as multitask regression of communicative, semantic, and syntactic skills. We found that domain adaptation before finetuning as well as question-answer pooling increased performance for the standard but not for the specialised models on all classification and regression tasks.

## 1 Introduction

Aphasia is a language impairment due to brain damage, after a stroke, traumatic head injury, brain tumours, or progressive neurological conditions. Depending on the brain regions affected, aphasia is featured differently. The most common types of aphasia are: global, amnesic (anomic), Wernicke's and Broca's aphasia (Caplan, 2003; Ardila, 2010). In Broca's aphasia, patients typically exhibit phonemic substitutions and have a non-fluent speech pattern. Wernicke's aphasia is characterised by an effortless but nonsensical speech. Global aphasia combines aspects of both Broca's and Wernicke's aphasia. Amnesic aphasia is primarily characterized by word retrieval and naming problems. Aphasia subtype classification is not straightforward and it is common that various aphasia types co-exist (Fridriksson et al., 2018).

Effective evidence-based therapy consists of high-intensity Speech-Language Therapy (SLT) which has been shown to improve linguistic capabilities (Peitz et al., 2024). However, this needs to be based on detailed diagnostics using appropriate tests. In German-speaking countries, the most common test used for aphasia diagnosis and monitoring is the standardised Aachen Aphasia Test (AAT) (Huber et al., 2013; Huber, 1983). This comprehensive test is designed to assess various aspects of language function, including comprehension, expression, repetition, and naming skills. It also provides information of probabilistic aphasia subtype and severity (Kohlschein et al., 2018). It consists of an examination of spontaneous language and five subtests. A 10-minute semi-structured interview, recorded during therapy, is rated in six domains: communicative behaviour, articulation/prosody, automatised language, semantics, phonology and syntax (Kohlschein et al., 2017). However, AAT is time-consuming and its result depends highly on the rater (Kohlschein et al., 2018), which usually is a highly trained speech and language therapist. An automatic aphasia diagnosis based on the AAT could help reduce waiting periods for patients and clinicians' burden as well as provide personalised remote rehabilitation strategies.

Prior work employing Machine Learning (ML) methods has explored aphasia and its subtype classification using connected speech, derived either from manual transcripts or Automatic Speech Recognition (ASR) systems (Fromm et al., 2022). These studies have focused on feature-based supervised methods, including traditional discourse features (e.g., syntactic complexity, proportion of nouns, verbs, adjectives) or embeddings by end-to-end approaches using large pre-trained models. Zusag *et al.* reported an F1 score of 0.84 for detecting amnesic aphasia, 0.77 for identifying Broca aphasia; and 0.78 for Wernicke aphasia using a Support Vector Classifier (SVC) and linguistic features (Zusag et al., 2023). Dunfield *et al.* employed sentence representation similarity features to capture symptoms of fluent aphasia and found a correlation of 0.61 with the Western Aphasia Battery-Revised Aphasia Quotient (Dunfield and Neumann, 2020). These features include question-answer similarity, closest question-answer pair identification, and binary sentence pair classification. The latter was obtained using BERT to predict the likelihood of a given sentence pair being related (Dunfield and

Neumann, 2020). Cong *et al.* leveraged Large Language Model (LLM)-surprisals to predict aphasia, its subtypes, and the level of severity. They reported an F1 score of 0.92 for predicting aphasia from healthy controls and 0.79 F1 score for identifying aphasia subtypes (Cong et al., 2024b). In another work, Cong *et al.* further employed surprisal values of LLMs, including GPT-2, Llama2, and BERT, alongside utterance length, to predict aphasia and its subtypes. Their results demonstrated an F1-score of 0.61 for detecting aphasia and 0.86 for classifying its subtypes in Chinese. For English, they reported an F1-score of 0.79 for identifying aphasia and 0.54 for distinguishing its subtypes (Cong et al., 2024a).

The contributions of our work of automatised aphasia assessment are as follows: (1) Aphasia transcripts are atypical on the lexical, syntactic, and semantic level. Such transcripts are usually not contained in the training material of pre-trained models, which might lower their general applicability on such clinical data. We are going to address this potential shortcoming by domain adaptation as described in section 3.2. (2) Relevant information is expected not to be contained only in the patients' answers in isolation but also within the context of the underlying question. We are going to address this contextualization by embedding pooling alternatives as presented in section 3.3.

## 2 Data

The German dataset was collected within the autoAAT BMBF project. It contains spontaneous speech samples, manual transcripts, and their associated clinical scores from the AAT. Transcripts were anonymised by removing all personal information. This dataset is built on the work presented in (Kohlschein et al., 2018). Many patients provided more than one recording due to repeated treatment cycles. The scores comprise the aphasia type classification and linguistic skills assessment. Aphasia type is categorised into the four classes Amnesic, Broca, Global, and Wernicke; since the project focus is to automatise aphasia diagnosis for tailored SLT, the dataset does not contain a control group. Other types of aphasia, such as primary progressive aphasia or unclassifiable, have been excluded of the analysis due to data sparsity. Linguistic skills are assessed separately in various impairment levels and on an expert-annotated six point scale (with 0 being the most severe and 5

meaning no impairment). This study focuses on three linguistic impairment levels: communicative behaviour (understanding and responding to questions), semantic structure (word finding difficulties and semantic paraphasias), and syntactic structure (sentence completeness and complexity).

The dataset comprises 331 participants, 92 female, 239 male, with a mean age of $53 \pm 13$ years. The major aphasia types are represented by the following numbers: 105 Global, 70 Broca, 32 Wernicke, and 34 Amnesic. The rest of the participants correspond to the excluded classes. Due to data protection regulations, the dataset cannot be shared. The dataset was split into speaker-disjunct training, development (10%), and test (20%) sets stratified on the aphasia type of each speaker by means of *splitutils* (Reichel, 2024). A random seed of 42 was applied to ensure reproducibility. Texts were cleaned by removing transcriber comments and special annotation symbols. The linguistics skills scales ranging from 0 to 5 were re-scaled to the range $[0, 1]$.

## 3 Methods

### 3.1 Modelling variants

For both tasks, aphasia type classification and linguistic skills regression, we started from two different base models: the general-purpose model *dbmdzbert-base-german-uncased* (Devlin et al., 2019) (referred to as *standard encoder* in the following), and *GerMedBERT/medbert-512* (Bressem et al., 2023), which was pre-trained on medical documents for applications in the clinical domain, henceforth referred to as *specialised encoder*.

For each of these encoders, we further created a variant domain-adapted to our specific aphasia dataset as described in section 3.2. Each of these four variants we combined with three different pooling architectures as described in section 3.3. We finetuned each of these 12 model variants on the two clinical tasks with 5 different random seeds, which we describe in section 3.4.

### 3.2 Domain adaptation

For domain adaptation, we followed the recipe of (Lendvai et al., 2023) applying vocabulary extension and Masked Language Modelling (MLM). We applied a 90/10 speaker disjunct and aphasia-label stratified split of the training partition into MLM training and development partition. Based on the MLM training partition we extended the tokeniz-

ers' vocabularies with the lexical content of the transcripts by adding up to 300 most frequent, yet unknown words with a minimum length of five characters. Subsequently, each base model was finetuned on the MLM task with a standard BertForMaskedLM head. Finetuning was done in 20 epochs with the AdamW optimizer, a learning rate of $2e-5$, a perplexity loss, and a batch size of 16. We kept the best model in terms of the lowest loss for the development set.

### 3.3 Pooling

We applied three types of pooling of the last hidden states of the encoder:

**a**: answer-only; we extract the embeddings only for the patient's answer and apply mean pooling of these embeddings;

**qa-c**: answer contextualised by question; we concatenate question and answer with a [SEP] token as for text entailment tasks (Putra et al., 2024), extract the embeddings for this text pair, and apply mean pooling on the answer part of this pair only, which is forwarded to the classification head;

**qa-cc**: answer contextualised by question plus question-answer coherence; as for *qa-c* we concatenate question and answer. Then, we concatenate the initial CLS token embedding with the mean embedding of the answer. This concatenated pooling we forward to the classification head.

Schematically, the pooling variants can be expressed as follows (the underlined constituents go into the pooling):

|  |  |
|---|---|
| **a:** | [CLS] <u>**answer**</u> |
| **qa-c:** | [CLS] question [SEP] <u>**answer**</u> |
| **qa-cc:** | <u>[**CLS**]</u> question [SEP] <u>**answer**</u> |

We expect *qa-cc* to capture not only answer contextualisation but also question-answer coherence due to the 'semantics' of the CLS token. Since this token had been pre-trained on the next sentence prediction task, it is expected to represent the information the pre-[SEP] text part contains about the post-[SEP] text part, which can be considered as an aspect of text coherence.

In total, we get 12 model variants defined by all combinations of **encoder type** (*standard, specialized*), **domain adaptation** (*yes, no*) and **pooling** (*a, qa-c, qa-cc*). The finetuning of these models on the two downstream tasks is described in the subsequent section 3.4.

### 3.4 Finetuning

**Architecture:** To each encoder, we add a two-layer head with a non-linear (tanh) layer and a linear output projection. For classification, this output projection has 4 outputs, one per aphasia type. For multitask regression, it has 3 outputs, one for communicative, semantic, and syntactic skills, respectively.

**Hyperparameters:** Each model was finetuned in 8 epochs with the AdamW optimizer, a learning rate of $3e-5$ and an effective batch size of 32. For classification, we used the weighted cross entropy loss and unweighted average recall (UAR) as metrics to be maximised on the development set. For regression, we used a Concordance Correlation Coefficient (CCC) loss and CCC metrics for the development set. We kept the models performing best on the development set for further evaluation on the test partition. Finetuning and evaluation was repeated five times with different random seeds (1, 9, 20, 21, 42, generated with *numpy.random.default_rng()*).

## 4 Results

Figures 1 and 2 show the results in terms of UAR and mean CCC for aphasia type classification and linguistic skills regression, respectively. As an overall tendency for the standard encoder, we observe that domain adaptation as well as question-answer contextualisation slightly improve the performances for classification as well as for regression, but not so for the specialised encoder.

The best aphasia type classification result, a UAR of $0.653$ averaged over all random seeds, was obtained with the standard encoder, and the *qa-cc* pooling variant accounting for contextualisation and coherence. For linguistic skills multitask regression, again, the standard encoder this time with the *qa-c* pooling variant for contextualisation only performed best, yielding a mean CCC of $0.755$ averaged over all random seeds. Split into the linguistic dimensions it achieved a CCC of $0.738$ for communicative, $0.695$ for semantics, and $0.831$ for syntactic skills prediction.

## 5 Discussion and Conclusion

We identified two challenges for finetuning pre-trained transformer models with aphasia data: First, this text data is rather atypical and usually not part of pre-training datasets. This missing link was addressed by domain adaptation. Second, patients' answers are not only to be seen in isolation but also within context with the corresponding ques-
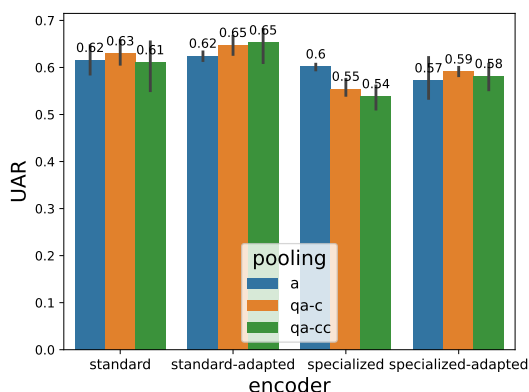
Figure 1: Aphasia type classification results: Un-weighted average recall (UAR) values for all encoder and pooling variant combinations (see section 3). Error bars indicate 95% confidence intervals.
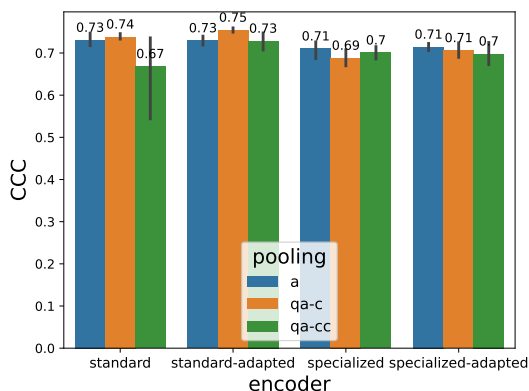


Figure 2: Communication, semantics, and syntactic skills multitask regression results: arithmetic mean Concordance correlation coefficient (CCC) values over the three regression dimensions for all encoder and pooling variant combinations (see section 3). Error bars indicate 95% confidence intervals.

tion. This contextualisation and coherence assessment was addressed by introducing different kinds of question-answer poolings.

For the standard encoder, domain adaptation as well as question-answer pooling turned out to be beneficial for both aphasia type classification as well as linguistic skills regression. Both strategies, by a low margin but consistently, lead to increased performance. As to pooling, for aphasia type classification, joint contextualisation and coherence assessment worked best, for regression contextualisation only lead to the highest performance.

The specialised encoder overall yielded lower performances compared to the standard encoder, which on first sight might appear counter-intuitive.

However, the specialised model was not necessarily expected to work better for patient data classification in the first place, since the pre-training material consists exclusively of expert texts from scientific publications and dictionaries, as reported in (Bressem et al., 2023). These documents usually do not include a large amount of patient transcripts, but rather few illustrative examples only. Therefore, this specialised model is well suited for tasks such as clinical expert text classification, but not necessarily for patient transcript classification. One major reason for the overall lower performance of the expert model might be that the specialised pre-training material contains much less variability than the standard encoder's pre-training data, so that it is less capable to extrapolate to that kind of data. Likely due to this shortcoming, the specialised model neither could profit from domain adaptation nor question-answer pooling.

For question-answer pooling, longer error bars were observed for *qa-cc* as opposed to *qa-c* in Figures 1 and 2. This indicates that joint contextualisation and coherence assessment is less stable across random seed variations than contextualisation alone, so that the latter seems to be preferable in terms of model robustness.

To conclude and to give an outlook, our results show that for the given data, aphasia modelling works best with domain-adapted standard BERT models with contextualised mean pooling of the embeddings of patients' utterances. These results were obtained on narrow manual transcripts that preserve linguistic peculiarities relevant for aphasia assessment. For a fully automated aphasia assessment, such transcripts would need to be generated by ASR models, that keep track of clinically relevant utterance characteristics such as disfluencies; see, e. g., (Zusag et al., 2023; Mihajlik et al., 2024; Gohider and Basir, 2024) for such ASR methods. Our next steps thus will include combining automated narrow transcription with our aphasia modelling approach.

**Ethics Statement**

This research was conducted with strict adherence to ethical standards. The dataset employed in this work was collected under the ethics approval number EK 23-125 by the Ethics Committee of the Medical Faculty of RWTH Aachen University. To further ensure privacy, audio and text data was anonymised removing all personal information.

Data was analysed transparently, avoiding bias and ensuring accuracy.

## Acknowledgements

## References

Alfredo Ardila. 2010. A proposed reinterpretation and reclassification of aphasic syndromes. *Aphasiology*, 24(3):363–394.

Keno K. Bressem, Jens-Michalis Papaioannou, Paul Grundmann, Florian Borchert, Lisa C. Adams, Leonhard Liu, Felix Busch, Lina Xu, Jan P. Loyen, Stefan M. Niehues, Moritz Augustin, Lennart Grosser, Marcus R. Makowski, Hugo JWL. Aerts, and Alexander Löser. 2023. MEDBERT.de: A comprehensive German BERT model for the medical domain. *arXiv preprint arXiv:2303.08179*. Keno K. Bressem and Jens-Michalis Papaioannou and Paul Grundmann contributed equally.

David Caplan. 2003. Aphasic syndromes. *Clinical neuropsychology*, 4:14–34.

Yan Cong, Jiyeon Lee, and Arianna LaCroix. 2024a. Leveraging pre-trained large language models for aphasia detection in English and Chinese speakers. In *Proceedings of the 6th Clinical Natural Language Processing Workshop*, pages 238–245, Mexico City, Mexico. Association for Computational Linguistics.

Yonghao Cong, Amy N. LaCroix, and Jiyeon Lee. 2024b. Clinical efficacy of pre-trained large language models through the lens of aphasia. *Scientific Reports*, 14:15573.

Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. 2019. Bert: Pre-training of deep bidirectional transformers for language understanding. In *Proceedings of the 2019 conference of the North American chapter of the association for computational linguistics: human language technologies, volume 1*, pages 4171–4186.

Katherine Dunfield and Günter Neumann. 2020. Automatic quantitative prediction of severity in fluent aphasia using sentence representation similarity. In *Proceedings of RaPID-2020 at LREC-2020*.

Julius Fridriksson, Dirk-Bart den Ouden, Argye E. Hillis, Gregory Hickok, Chris Rorden, Alexandra Basilakos, Grigori Yourganov, and Leonardo Bonilha. 2018. Anatomy of aphasia revisited. *Brain*, 141(3):848–862.

Davida Fromm, Joel Greenhouse, Molly Pudil, Yiwen Shi, and Brian MacWhinney. 2022. Enhancing the classification of aphasia: A statistical analysis using connected speech. *Aphasiology*, 36(12):1492–1519.

Nada Gohider and Otman A Basir. 2024. Recent advancements in automatic disordered speech recognition: A survey paper. *Natural Language Processing Journal*, 9:100110.

Walter Huber. 1983. *Aachener Aphasie Test (AAT)*. Verlag für Psychologie Dr. C.J. Hogrefe.

Walter Huber, Klaus Poeck, and Luise Springer. 2013. *Klinik und Rehabilitation der Aphasie: Eine Einführung für Therapeuten, Angehörige und Betroffene*. Georg Thieme Verlag.

Christian Kohlschein, Daniel Klischies, Björn Schuller, Tobias Meisen, and Cornelius Johannes Werner. 2018. Automatic processing of clinical aphasia data collected during diagnosis sessions: Challenges and prospects. In *International Conference on Language Resources and Evaluation*.

Christian. Kohlschein, Maximilian Schmitt, Björn Schuller, Sabina Jeschke, and Cornelius J. Werner. 2017. A machine learning based system for the automatic evaluation of aphasia speech. In *2017 IEEE 19th International Conference on e-Health Networking, Applications and Services (Healthcom)*. IEEE.

Piroska Lendvai, Uwe Reichel, Anna Jouravel, Achim Rabus, and Elena Renje. 2023. Domain-adapting BERT for attributing manuscript, century and region in pre-modern Slavic texts. In *Proceedings of the 4th Workshop on Computational Approaches to Historical Language Change*, pages 15–21.

Péter Mihajlik, Yan Meng, Máté S Kádár, Julian Linke, Barbara Schuppler, and Katalin Mády. 2024. On disfluency and non-lexical sound labeling for end-to-end automatic speech recognition. In *Interspeech 2024*, pages 1270–1274, Kos Island, Greece.

Dorothea Peitz, Beate Schumann-Werner, Katja Hussmann, Joao Pinho, Hong Chen, Ferdinand Binkofski, Walter Huber, Klaus Willmes, Stefan Heim, Jörg B. Schulz, Bruno Fimm, and Cornelius J. Werner. 2024. Success rates of intensive aphasia therapy: Real-world data from 448 patients between 2003 and 2020. *Journal of Neurology*.

I Made Suwija Putra, Daniel Siahaan, and Ahmad Saikhu. 2024. Recognizing textual entailment: A review of resources, approaches, applications, and challenges. *ICT Express*, 10(1):132–155.

Uwe Reichel. 2024. splitutils v0.3.0. Zenodo.

Markus Zusag, Lisa Wagner, and Thomas Bloder. 2023. Careful whisper – leveraging advances in automatic speech recognition for robust and interpretable aphasia subtype classification. In *Proceedings of Interspeech 2023*, pages 3013–3017.

# A Retail-Corpus for Aspect-Based Sentiment Analysis with Large Language Models

**Oleg Şilcenco[1], Marcos R. Machado[1], Wallace C. Ugulino[1], Daniel Braun[2]**
[1]University of Twente, [2]Marburg University

olegsilcenco@gmail.com, {m.r.machado,w.corbougulino}@utwente.nl, daniel.braun@uni-marburg.de

## Abstract

Aspect-based sentiment analysis enhances sentiment detection by associating it with specific aspects, offering deeper insights than traditional sentiment analysis. This study introduces a manually annotated dataset of 10,814 multilingual customer reviews covering brick-and-mortar retail stores, labeled with eight aspect categories and their sentiment. Using this dataset, the performance of GPT-4 and LLaMA-3 in aspect based sentiment analysis is evaluated to establish a baseline for the newly introduced data. The results show both models achieving over 85% accuracy, while GPT-4 outperforms LLaMA-3 overall with regard to all relevant metrics.

## 1 Introduction

Sentiment analysis, i.e. the automatic identification of the sentiment expressed in, e.g., a text, is a widely used technique in research, business, politics, and many other domains (Wankhade et al., 2022). While traditional sentiment analysis methods focus mainly on detecting sentiment at the sentence or document level, aspect-based sentiment analysis is a more fine-grained approach through which particular aspects expressed in a text are identified with their corresponding sentiment (Zhang et al., 2022). For a movie review, e.g., in this way not just the overall sentiment expressed by the review is identified but also, for example, whether the reviewer liked or disliked the score or camera work. Such more fine-grained insights are particularly valuable to businesses as they provide better insights into the needs of customers.

While datasets for traditional sentiment analysis are widely available (see e.g. Tan et al. (2023); Kenyon-Dean et al. (2018); Wagh and Punde (2018), and Saif et al. (2013) for an overview of popular datasets), in part because they can be gathered in an automated fashion from services that use a combination of a score (e.g. in the form

of stars) alongside with a textual review, the number of datasets available for aspect-based sentiment analysis is much more restricted (see e.g. Nazir et al. (2022) and Hua et al. (2024)). Additionally, most of the existing datasets either contain a small number of aspects per item or all aspects in one item have the same polarity (Jiang et al., 2019).

In this paper, we introduce a new, manually annotated, dataset for aspect-based sentiment analysis, that consists of 10,814 reviews for brick-and-mortar retail stores, scraped from Google Maps. The reviews cover different countries and languages and have been annotated with eight different aspect categories (see Table 3) resulting in a total of 16,994 labels. The dataset is available on GitHub[1]. A detailed datasheet (Gebru et al., 2021) for the corpus can be found in Appendix B.

In addition, we present a Large Language Model (LLM)-based baseline for the newly introduced dataset comparing the performance of Meta's LLaMa-3 and OpenAI's GPT-4. The results show that while both models perform well with an accuracy of more than 85%, GPT-4 consistently outperforms LLaMa-3 across aspects and metrics.

## 2 Related Work

Compared to "traditional" sentiment analysis, aspects-based sentiment analysis presents a more complex challenge, encompassing two distinct stages: identifying all aspects described, and subsequently determining the sentiment towards each aspect.

### 2.1 LLMs for Aspect-Based Sentiment Analysis

As for most NLP applications, recent literature about aspect-based sentiment analysis has mainly focused on the performance of LLMs. Recent

---

[1] https://github.com/Responsible-NLP/ABSA-Retail-Corpus

Table 1: Sample of scraped data

| Country | City | Published At | Text | Stars |
|---------|------|--------------|------|-------|
| Belgium | Maasmechelen | 18-03-2023 | Najbolja i najkvalitetnija roba | 5 |
| France | Serris | 29-12-2019 | Great prices | 5 |
| Italy | Marcianise | 21-11-2023 | Ho scoperto questo negozio grazie... | 5 |
| Belgium | Mechelen | 30-11-2017 | Mooie propere zaak maar verkoper... | 3 |

studies that compared the performance of LLMs against smaller language models (SLMs), like BERT, across a spectrum of sentiment analysis tasks, including conventional sentiment classification and aspect-based analysis, found that LLMs exhibit proficiency in simpler tasks, such as sentiment classification, but encounter difficulties in tasks requiring nuanced understanding or structured sentiment information (Zhang et al., 2024; Macháliková, 2023; Han and Moghaddam, 2023).

In few-shot learning scenarios, however, where annotation resources are limited, LLMs have shown superior performance (Zhang et al., 2024). According to Magdaleno et al. (2024), LLMs also outperform smaller models in such tasks as predicting ratings of businesses based on online reviews, and leveraging aspect-based sentiment analysis techniques. Moreover, LLMs have introduced innovative methodologies for context-aware analysis, as shown by Jeong and Lee (2024) in the context of hotel complaint reviews.

Comparative analyses between GPT-3.5, BERT, RoBERTa, and LLaMA report superior performance of GPT-3.5, specifically in predicting product review ratings post fine-tuning (Roumeliotis et al., 2024).

Krugmann and Hartmann (2024) report that using a zero-shot nature, LLMs can not only compete with but in some cases also surpass traditional transfer learning methods in terms of sentiment classification accuracy. Additionally, studies emphasize the competitive performance of GPT-3.5 model in discerning nuanced sentiments like irony within social media tweets, achieved through *prompt engineering* without explicit training (Carneros-Prado et al., 2023) Moreover, effective prompting engineering and *fine-tuning* are identified as crucial factors for achieving enhanced outputs and cost efficiency, further accentuating the potential of LLMs in customer satisfaction analysis and industry practices (Roumeliotis et al., 2024).

Comparative studies between LLMs and lexicon-based methods, show that LLMs clearly outperform such methods, while being particularly good in annotating sentiment analysis data and achieving over 94% accuracy in long-form sentiment reviews from Twitter social media users and Amazon customers, owing to its prowess in handling emojis, sarcasm, and contextual nuances (Belal et al., 2023). Notably, the literature suggests that GPT's integration into business customer sentiment analysis reveals its potential to significantly enhance understanding of customer sentiments, offering valuable insights for decision-making processes by comprehending both general sentiments and nuanced factors within customer texts (Sudirjo et al., 2023).

## 2.2 Datasets for Aspect-Based Sentiment Analysis

Table 2 shows a list of existing datasets for aspect-based sentiment analysis. While a number of datasets exists, the vast majority of them only covers the English language. Multilingual data sets, such as the one introduced in this paper, are particularly rare. With more than 10,000 annotated instances, the dataset introduced in this article is also one of the largest data sets for aspect-based sentiment analysis that is currently available.

## 3 Corpus

The data collection for the corpus relied on Google Maps reviews due to their accessibility and richness. While other publicly available datasets, such as those from product reviews or social media platforms, exist, many lack the level of granularity necessary for aspect-based sentiment analysis. These datasets typically emphasize general sentiment or aggregate ratings, which limits their suitability for examining specific aspects of customer feedback. Google Maps represent a robust platform where customers can leave reviews about a specific store or location. With its vast user base and widespread popularity, Google Maps is one of the most prominent platforms for user-generated reviews.

To efficiently scrape the reviews from Google

Table 2: Existing datasets for aspect-based sentiment analysis

| Dataset | Domain | Lang. | Size | Sources |
|---------|--------|-------|------|---------|
| SemEval 2014 | Service and Product Reviews | English | 7,686 | Pontiki et al. (2014) |
| SemEval 2015 | Service and Product Reviews | English | 5,596 | Pontiki et al. (2015) |
| SemEval 2016 | Service and Product Reviews | 8 | 6,243 | Pontiki et al. (2016) |
| Pars-ABSA | Service Reviews | Persian | 5,602 | Shangipour ataei et al. (2022) |
| Foursqaure | Service Reviews | English | 585 | Brun and Nikoulina (2018) |
| ACOS | Service and Product Reviews | English | 6,362 | Cai et al. (2021) |
| SentiHood | Neighbourhood Q&A | English | 5,215 | Saeidi et al. (2016) |
| MAMS | Service Reviews | English | 13,854 | Jiang et al. (2019) |
| *Our dataset* | Service Reviews | 45 | 10,814 | |

Maps, the service Apify[2] was utilized. Privacy and personal data were primary concerns during the data collection process. Despite Google Maps reviews being publicly accessible, it was important to be cautious to ensure compliance with privacy standards. Apify's functionality enabled the selection of only the essential columns, omitting personal identifiers such as the name/nickname of the reviewer. Consequently, the collected data only included the review text, star rating, timestamp, and the location of the store (country and city). A small sample of such can be seen in Table 1.

## 3.1 Data Cleaning and Augmentation

A total of 24,361 reviews were collected in that way. Many reviews contained only a star rating without any textual review, these entries were excluded from the dataset, resulting in a final dataset of 10,814 reviews. Since the dataset contains reviews in a variety of languages, an additional column was added to encode the language of each review. The Google Translate API, accessed via the Python library `googletrans`, was utilized to detect the language of the reviews. The languages are encoded in ISO-639 format. The dataset also contains entries with unidentified languages and entries that contain only symbols or emojis. Additionally, the publication time, which was in textual format in the raw data, was converted to ISO-8601 format.

## 3.2 Data Annotation

The most important decision that had to be made before the data annotation is the definition of aspect categories. *Aspect categories* are higher-level concepts that pool different *aspect terms* to allow

for more structured insights (Hua et al., 2024). The selection of the aspect categories was based on existing literature (particularly Kang et al. (2022); Ramaswami and Varghese (2003); Fakhira and Simanjuntak (2023)) and interviews with domain experts, to ensure that the chosen categories are relevant from both a scientific and a practitioner perspective. The eight aspect categories that have been derived from this process are shown in Table 3. They cover a broad range of customer experiences and provide valuable insights into different facets of the reviews and the businesses behind them. Each of these categories, if identified in a review, was assigned a sentiment label (negative, positive, or neutral).

The dataset was manually labeled by the authors. Manual labeling, though time-consuming, is critical for ensuring high-quality data annotations. Given the labor-intensive nature of this task, a custom labeling tool was developed to facilitate the process. Similar approaches have been employed in other studies. For instance, the SemEval-2014 task 4 involved creating annotation guidelines and tools for manual labeling to build benchmark datasets (Pontiki et al., 2014). Do et al. (2020) surveyed various tools and methods developed to assist in aspect-level sentiment annotation, while Li et al. (2012) presented a tool designed to create high-quality training data through manual annotation.

Figure 1 shows the tool that was developed to annotate the dataset. It allowed the authors to review the text, select the relevant aspect, and assign the corresponding sentiment. The tool features a user-friendly interface with functionalities such as language detection and translation to English, sentiment selection, and annotation saving.

10% of the data was annotated by two annotators to conduct an inter-annotator agreement study. Using Krippendorff's Alpha, the inter-annotator agree-

---

[2] https://apify.com/compass/google-maps-reviews-scraper

184

Table 3: Chosen aspect categories and their description

| Aspect | Explanation |
|---|---|
| *Product* | Encompasses clothing collection, item quality, variety, display, and selection. |
| *Service* | Includes staff, assistance, crew, employee attitude, handling, and hospitality. |
| *Brand* | Pertains to overall brand perception. |
| *Price* | Relates to the cost of products and services, including promotions and discounts. |
| *Store* | Covers specific shop location, atmosphere, and environment. |
| *Online* | Concerns the online ordering experience. |
| *Return* | Includes the experience of returning an item for both physical or online procurement. |
| *General* | Overall experience of the customer without a specific aspect mention. |



Figure 1: Labeling tool

ment is $\alpha = 0.71$, a value that is comparable to other ABSA datasets and in general indicates a reliable annotation. An in-depth analysis revealed only two instances in which the annotators chose the same aspect but different polarities. In both cases it was not a direct contradiction, as in a positive and a negative label at the same time, but one label was neutral. An analysis on aspect-level shows that more specific aspects like price and service show higher agreement, while more general categories like store and "general" show lower agreement. A reoccurring pattern is that often, annotators agree on most aspects, but one annotator adds a single additional aspect (like general), thereby decreasing agreement.

### 3.3 Data Analysis

Out of the 10,814 reviews in the dataset, 4,838 (or 44.7%) contain more than one aspect. On average, each review contains 1.6 aspects. Figure 2 shows how often each aspect category occurs in the dataset. The most frequently occurring aspect category is service, which is mentioned in 6,065 reviews. The least mentioned category, online, only occurred in 51 reviews, which is not surprising, given that the reviews were specifically collected for physical store locations.

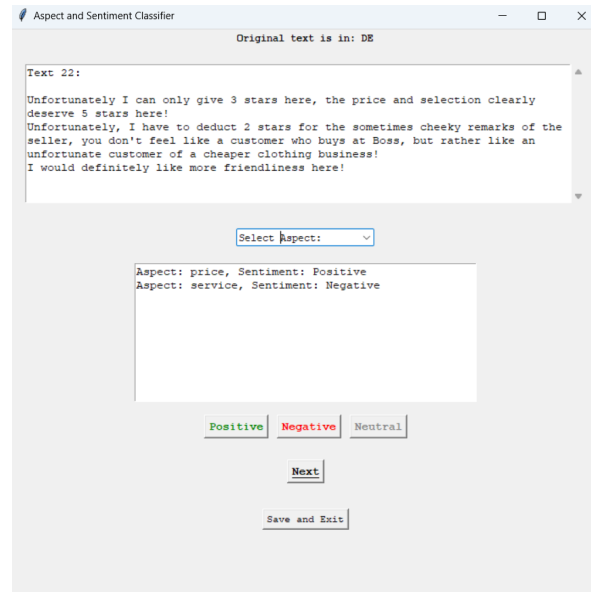The majority of the aspects mentioned in the dataset are positively connoted, as shown in Figure 3. Only for the aspect category return, the majority of the reviews expresses a negative sentiment. The most balanced aspect category is online, in which 55% of all mentioned aspects are positive, 8% neutral, and 37% negative. Neutral aspects are rarely mentioned across all categories. The highest share of neutral aspects can be found for the category general, with a little over 8%.

The reviews in the dataset have an average length of 121 characters, ranging from just one character (mostly emojis) to 3,735 characters. They cover stores from nine different European countries (Germany, France, Netherlands, Italy, Spain, Austria, Belgium, Portugal, and Switzerland; in descending frequency of occurrence) and are written in 45 different languages (see Figure 4 for a distribution of the languages).

## 4 Experimental Set-Up

To establish a baseline in aspect-based sentiment analysis for the newly introduced dataset, we conducted an experiment comparing the performance of the open weights model LLaMa-3 (Dubey et al., 2024) and the proprietary GPT-4 model (Achiam et al., 2023).

LLaMa-3 (Meta-Llama-3-70B-Instruct) was integrated through the HuggingFace library. To enhance performance and efficiency, quantization techniques were applied, reducing the model's weight precision via BitsAndBytesConfig, which optimized computational resources, particularly for
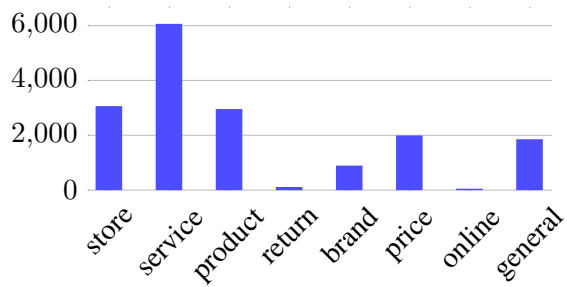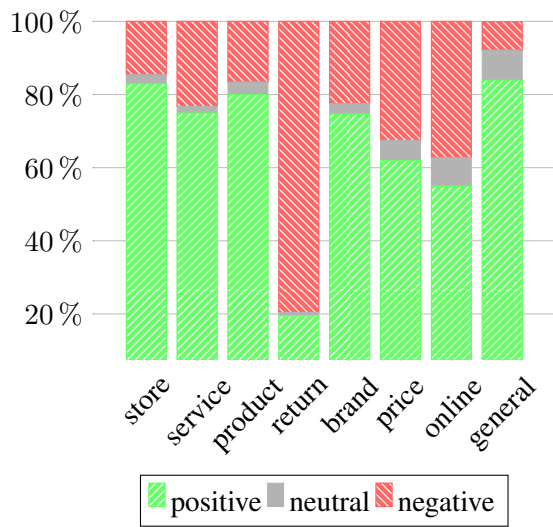
Figure 2: Occurrences of each aspect category



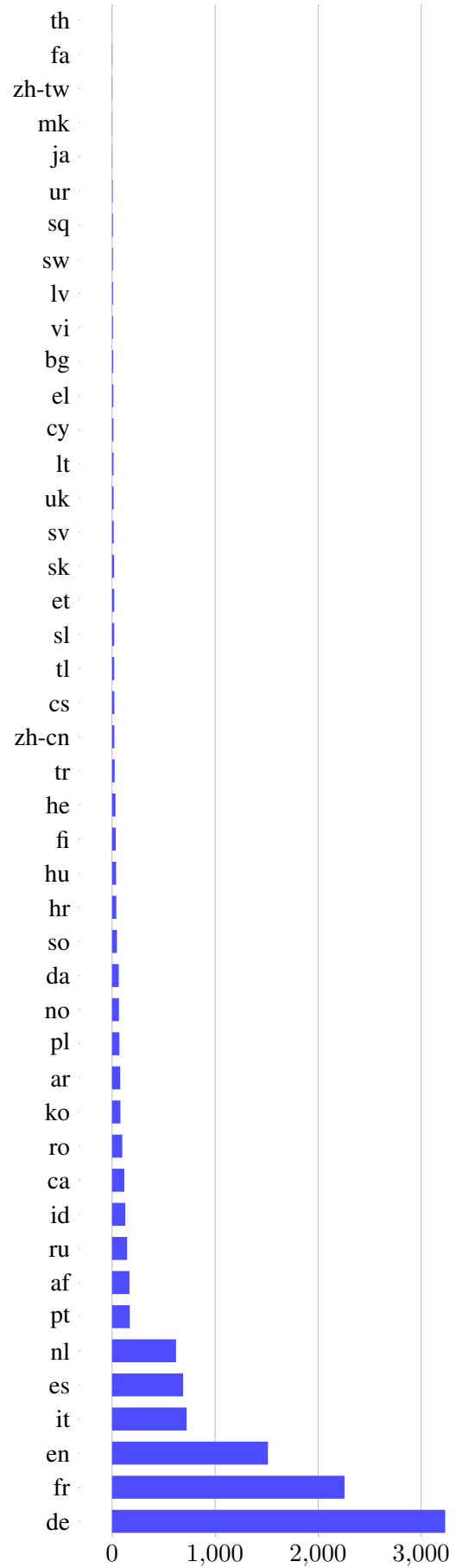Figure 3: Share of positive, neutral, and negative aspects per category



Figure 4: Number of reviews per language (ISO-639)

186

local deployment on GPU clusters.

GPT-4 was used through Azure's OpenAI Service. The implementation leveraged LangChain and its `PromptTemplate` component.

For both models, we used prompt engineering (Brown et al., 2020; Radford et al., 2019; Gao et al., 2021; Lu et al., 2021), and both system and user prompts, to facilitate the aspect-based sentiment analysis. *System prompts* are designed to define the role of the language model and establish operational guidelines to ensure consistency in responses. For both GPT-4 and LLaMA-3, the system prompt instructed the model to perform aspect-based sentiment analysis for a specific company. This approach is known as intent classification. The system prompt was implemented differently for each model, reflecting their respective frameworks. GPT-4 used LangChain's `SystemMessage` object to deliver the system instructions, while LLaMA-3 structured the system message as part of its chat template.

Despite syntactic differences, both implementations enforce a structured response format by defining system behavior upfront. LLaMA-3 specifies message roles such as *"system"* and *"user"* within its message list (see Listing 1), while GPT-4 utilizes LangChain's *langchain_core.messages* framework to differentiate between system and human messages (see Listing 2). These system prompts establish a consistent operational framework, ensuring the model generates precise and task-specific sentiment analysis responses.

User prompts provide the actual reviews and define the task parameters, guiding the model in identifying relevant aspects and classifying sentiments. The construction of these prompts is essential for ensuring accurate analysis, as they help the model distinguish between various aspects of a review and interpret sentiments effectively. To enhance performance, the prompts can incorporate structured instructions, few-shot examples, delimiters, and attention mechanisms.

One effective prompt engineering strategy is task decomposition, where a complex task is divided into smaller, more manageable steps. For instance:

> "First, identify the aspects in the provided review from the given list, and then find the customer sentiment (positive, neutral, or negative) for each of the aspects. ..."

Another key technique is few-shot prompting,

where examples of correct outputs are included in the prompt to guide the model's response.

> " ... You can follow the examples below:
>
> Review: The product quality is great but the customer service is terrible. Aspect: product Sentiment: positive Aspect: service Sentiment: negative
>
> Review: I love the location of the store. The collection and selection look great, however, the prices are too high. Aspect: store Sentiment: positive Aspect: product Sentiment: positive Aspect: prices Sentiment: negative
>
> Review: All top, everything as I expected, recommend. Aspect: general Sentiment: positive ... "

Finally, attention mechanisms can be influenced by placing key instructions at the beginning and end of the prompt. For example:

> "First, identify the aspects in the provided review from the given list, and then find the customer sentiment (positive, neutral, or negative) for each aspect.
>
> Make sure to take into account the difference in language, cultural aspects, sarcasm, emojis, and other linguistic behaviors when interpreting and assessing the reviews. ...
>
> ... Remember to strictly focus only on the aspects from the list and reply only with the answer in the following JSON format:
>
> ["aspect1": "sentiment", ... "aspectN": "sentiment"]"

By placing important contextual elements at the beginning and specifying output format at the end, the model is guided to prioritize crucial information while maintaining structured responses. The final prompts used in the experiment can be found in Appendix A.

## 5 Results

Table 4 shows the evaluation of the aspect-based sentiment analysis experiment. Overall, GPT-4 outperformed LLaMA-3 in every single metric, with the widest gap in precision and the narrowest gap in recall. Given the imbalance between positive

Table 4: Precision, Recall, F1 Score, and Accuracy of the aspect-based sentiment analysis per model and aspect

| Aspect | GPT-4 | | | | LLaMA-3 | | | |
|---|---|---|---|---|---|---|---|---|
| | Prec. | Recall | F1 | Accur. | Prec. | Recall | F1 | Accur. |
| Store | **73.00%** | 73.61% | **73.31%** | 71.23% | 59.02% | **80.13%** | 67.97% | **75.53%** |
| Service | **93.88%** | **98.00%** | **95.89%** | **95.31%** | 93.43% | 96.44% | 94.91% | 94.05% |
| Product | **65.34%** | 92.98% | **76.75%** | **88.39%** | 58.12% | **93.70%** | 71.74% | 87.82% |
| Return | 55.13% | **77.48%** | **64.42%** | **76.11%** | **62.89%** | 54.46% | 58.37% | 53.98% |
| Brand | **57.92%** | 70.65% | **63.66%** | 65.66% | 39.81% | **84.92%** | 54.21% | **80.43%** |
| Price | **82.71%** | **85.86%** | **84.26%** | **79.23%** | 79.04% | 80.84% | 79.93% | 73.09% |
| Online | 29.10% | **82.98%** | **43.09%** | **76.47%** | **30.09%** | 72.34% | 42.50% | 66.67% |
| General | **49.75%** | **78.09%** | **60.78%** | **73.92%** | 43.96% | 65.13% | 52.49% | 62.13% |
| **Micro avg.** | **74.48%** | **87.58%** | **80.50%** | **83.80%** | 66.84% | 86.88% | 75.55% | 82.57% |
| **Macro avg.** | **63.35%** | **82.46%** | **70.27%** | **78.29%** | 58.30% | 78.49% | 65.27% | 74.21% |

and negative reviews, the difference in accuracy of both models is only about 1.2 percentage points, despite the larger difference in precision. Only for the aspect categories "store" and "brand", LLaMa-3 outperformed GPT-4 in respect to accuracy.

In high-frequency aspects such as "service" and "product", both models showed strong performance, with GPT-4 recording accuracies of 95.31% and 88.39%, respectively, compared to LLaMA-3's 93.05% and 87.82%. Notably, LLaMA-3 demonstrated improvements in the "product" category, narrowing the gap with GPT-4. These results indicate that both models reliably handle frequently mentioned aspects, although GPT-4 retains a slight edge in overall robustness.

When addressing lower-frequency or more nuanced aspects such as "return", "brand", and "online", both models continued to face challenges, albeit with notable differences. GPT-4 demonstrated better performance in the "return" category, achieving an accuracy of 76.11% compared to LLaMA-3's 53.98%. Similarly, in "online", GPT-4 outperformed LLaMA-3, recording accuracies of 76.47% versus 66.67%. These findings underscore GPT-4's greater capability to handle complex sentiment categories, though significant gaps remain. Both models struggled with the "online" aspect in terms of precision, with GPT-4 achieving 29.10% and LLaMA-3 slightly higher at 30.09%. These metrics highlight a broader limitation in capturing context-dependent nuances in less straightforward sentiment categories. In order for an aspect-based sentiment to be classified as correct, both the aspect and the sentiment expressed with it have to be extracted correctly. Notably, a large share of errors already occurs during the identification of aspects

(see Table 5). The number of aspects that have been identified correctly but the sentiment was misclassified is relatively small.

# 6 Error Analysis

A deeper analysis of the errors made by both models revealed that LLaMA-3 exhibited a tendency to incorrectly identify aspects that are not present in the text based on mentioned keywords. For instance, in reviews such as "Great store for its amazing service and help from the assistants", LLaMA-3 frequently identified "store" as an aspect, whereas the correct interpretation according to the aspect categories defined in Table 3 would be "service". This over-sensitivity to mentions of keywords is also visible in Table 4, where the precision for the aspect store is particularly low for LLaMA-3. Yet, given the prevalence of the aspect across the dataset, this over-sensitivity might also partially explain why this is one of just two aspect categories in which LLaMA-3 achieved a higher recall and accuracy than GPT-4.

Both models encountered difficulties in consistently handling aspects from the categories "general" and "brand". Reviews with broad or ambiguous sentiments about a brand in general, such as "Brand for Bosses," or "I love BOSS" for the clothing brand "BOSS" posed a significant challenge for both models. These reviews often led to inconsistent labeling, with models sometimes assigning both "brand" and "general" aspects or failing to distinguish between them altogether. Given the ambiguity of such statements, even human annotators would likely struggle to reach consensus, making this a particularly challenging area for automated analysis and labeling.

Table 5: Aspect identification accuracy

| Aspect | GPT-4 | LLaMA-3 |
|---|---|---|
| Store | 74.47% | **81.26%** |
| Service | **98.05%** | 96.57% |
| Product | 93.33% | **94.10%** |
| Return | **77.88%** | 54.87% |
| Brand | 72.73% | **85.71%** |
| Price | **86.95%** | 82.69% |
| Online | **84.31%** | 74.51% |
| General | **79.27%** | 66.74% |
| **Micro avg.** | **88.12%** | 87.57% |
| **Macro avg.** | **83.37%** | 79.44% |

LLaMA-3 showed a tendency to classify reviews under the apsect category "brand" more frequently than GPT-4, which contributed to its lower accuracy for the aspect category "general", scoring 66.74% comparing to 79.27% for GPT-4. Simultaneously, LLaMa-3 outperformed GPT-4 in the aspect category "brand", being one of only two aspect categories where it outperformed GPT-4 model. This behavior, while indicative of an attempt to capture a broader range of sentiment, often led to increased misclassification rates when the sentiment was intended to be more general. The higher labeling frequency for "brand" by LLaMA-3 also aligns with its overall lower precision and F1 scores in this aspect, reflecting a trade-off between recall and precision.

## 7 Conclusion

This paper introduced a new multilingual corpus for aspect-based sentiment analysis that is based on more than 10,000 reviews of brick-and-mortar stores and was manually labeled with eight aspect categories, namely product, service, brand, price, store, online, return, and general (see Table 3).

Additionally, an experiment was conducted to establish a baseline for LLM-based aspect-based sentiment analysis on the newly introduced corpus by comparing the performance of GPT-4 and LLaMA-3. The results indicate that both models proficiently identify elements and attitudes in customer reviews, with GPT-4 continuously surpassing LLaMA-3 in precision, recall, and accuracy. Although both models achieved accuracy exceeding 85%, they performed insufficiently for the "store" and "brand" aspects, indicating areas for enhancement. From an application perspective, in addition to the performance, it is also worth considering the

potential costs of different approaches and models. At the time of conducting the experiments, GPT-4 was queried through the Azure OpenAI API at a total cost of $240.60 for the complete dataset, or an average cost of $0.022 per review. With the price for one million input tokens being around $2.50 and the price for one million output tokens being around $10.00. LLaMA-3, on the other hand, cannot just be self hosted, but also used through cloud providers like groq who charge in the realm of $0.59 / $0.79 per one million input / output tokens, providing significantly cheaper options.

## Ethics

By using public reviews and ensuring during both collection and annotation of the data that no identifiable information is contained in the reviews, we tried to ensure that our work has no direct adverse effects to anyone. Nevertheless, given the nature of the task, companies could use an approach like the one outlined in this work to try to automatically assess job performance of workers in physical store locations (e.g. by focusing on the aspect category service). However, we believe that in practice, the danger for such applications is relatively low given that most reviews (and none in our dataset) name specific employees or provide other information that could be used for the identification of individual employees like exact data and time of an interaction. While the environmental impact of LLMs is mostly discussed with regard to their training and fine-tuning, ever larger models also have an increasingly significant environmental impact during inference, which also holds true for this work.

## Limitations

This study faced limitations inherent to the rapidly evolving field of LLM research. The introduction of newer models may render some aspects of this study's model selection less timely, although its fundamental methodologies remain relevant. Computational limitations also restricted fine-tuning and advanced quick engineering for LLaMA-3, potentially affecting its performance. Furthermore, the hand annotated dataset, albeit comprehensive, introduced subjectivity in aspect and sentiment classification, with ambiguous phrases and linguistic variances presenting hurdles to consistency. Lastly, the predefined aspect categories may not generalize across all use cases, making prompt design sensitive to initial definitions.

# References

Josh Achiam, Steven Adler, Sandhini Agarwal, Lama Ahmad, Ilge Akkaya, Florencia Leoni Aleman, Diogo Almeida, Janko Altenschmidt, Sam Altman, Shyamal Anadkat, et al. 2023. Gpt-4 technical report. *arXiv preprint arXiv:2303.08774*.

Mohammad Belal, James She, and Simon Wong. 2023. Leveraging chatgpt as text annotation tool for sentiment analysis. *arXiv preprint arXiv:2306.17177*.

Tom Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, et al. 2020. Language models are few-shot learners. *Advances in neural information processing systems*, 33:1877–1901.

Caroline Brun and Vassilina Nikoulina. 2018. Aspect based sentiment analysis into the wild. In *Proceedings of the 9th Workshop on Computational Approaches to Subjectivity, Sentiment and Social Media Analysis*, pages 116–122, Brussels, Belgium. Association for Computational Linguistics.

Hongjie Cai, Rui Xia, and Jianfei Yu. 2021. Aspect-category-opinion-sentiment quadruple extraction with implicit aspects and opinions. In *Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing (Volume 1: Long Papers)*, pages 340–350, Online. Association for Computational Linguistics.

David Carneros-Prado, Laura Villa, Esperanza Johnson, Cosmin C Dobrescu, Alfonso Barragán, and Beatriz García-Martínez. 2023. Comparative study of large language models as emotion and sentiment analysis systems: A case-specific analysis of gpt vs. ibm watson. In *International Conference on Ubiquitous Computing and Ambient Intelligence*, pages 229–239. Springer.

Bao Lieu Do, Lam Pham Huy, and Xuan-Linh Tran. 2020. Automated tools for aspect-level sentiment analysis: A survey. *Information Processing & Management*, 57(6):102311.

Abhimanyu Dubey, Abhinav Jauhri, Abhinav Pandey, Abhishek Kadian, Ahmad Al-Dahle, Aiesha Letman, Akhil Mathur, Alan Schelten, Amy Yang, Angela Fan, et al. 2024. The llama 3 herd of models. *arXiv preprint arXiv:2407.21783*.

Nisrina Nur Fakhira and Megawati Simanjuntak. 2023. Content analysis of consumer reviews and comments on e-commerce. *Jurnal Doktor Manajemen*, 6:2.

Tianyu Gao, Adam Fisch, and Danqi Chen. 2021. Making pre-trained language models better few-shot learners. *arXiv preprint arXiv:2012.15723*.

Timnit Gebru, Jamie Morgenstern, Briana Vecchione, Jennifer Wortman Vaughan, Hanna Wallach, Hal Daumé III au2, and Kate Crawford. 2021. Datasheets for datasets. *Preprint*, arXiv:1803.09010.

Yi Han and Mohsen Moghaddam. 2023. A design knowledge guided position encoding methodology for implicit need identification from user reviews. In *International Design Engineering Technical Conferences and Computers and Information in Engineering Conference*, volume 87295, page V002T02A092. American Society of Mechanical Engineers.

Yan Cathy Hua, Paul Denny, Jörg Wicker, and Katerina Taskova. 2024. A systematic review of aspect-based sentiment analysis: domains, methods, and trends. *Artificial Intelligence Review*, 57(11):296.

Nayoung Jeong and Jihwan Lee. 2024. An aspect-based review analysis using chatgpt for the exploration of hotel service failures. *Sustainability*, 16(4):1640.

Qingnan Jiang, Lei Chen, Ruifeng Xu, Xiang Ao, and Min Yang. 2019. A challenge dataset and effective models for aspect-based sentiment analysis. In *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)*, pages 6280–6285, Hong Kong, China. Association for Computational Linguistics.

Min Kang, Bing Sun, Tian Liang, and Hong-Ying Mao. 2022. A study on the influence of online reviews of new products on consumers' purchase decisions: An empirical study on jd. com. *Frontiers in Psychology*, 13:983060.

Kian Kenyon-Dean, Eisha Ahmed, Scott Fujimoto, Jeremy Georges-Filteau, Christopher Glasz, Barleen Kaur, Auguste Lalande, Shruti Bhanderi, Robert Belfer, Nirmal Kanagasabai, Roman Sarrazingendron, Rohit Verma, and Derek Ruths. 2018. Sentiment analysis: It's complicated! In *Proceedings of the 2018 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long Papers)*, pages 1886–1895, New Orleans, Louisiana. Association for Computational Linguistics.

Jan Ole Krugmann and Jochen Hartmann. 2024. Sentiment analysis in the age of generative ai. *Customer Needs and Solutions*, 11(1):3.

Fang Li, Yulan He, Weizhi Meng, and Kun Wu. 2012. A generic approach for sentiment analysis and opinion mining. *International Journal of Computer Applications*, 49(3):21–28.

Xinxi Lu, Antoine Bosselut, Junyi Jessy Hao, and Yejin Choi. 2021. Fantastically ordered prompts and where to find them: Overcoming few-shot prompt order sensitivity. In *Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing (Volume 1: Long Papers)*, pages 4074–4088.

Kristina Macháliková. 2023. Utilizing chatgpt for sentiment analysis.

Diego Magdaleno, Martin Montes, Blanca Estrada, and Alberto Ochoa-Zezzatti. 2024. A gpt-based approach for sentiment analysis and bakery rating prediction. In *Advances in Computational Intelligence. MICAI 2023 International Workshops*, pages 61–76, Cham. Springer Nature Switzerland.

Ambreen Nazir, Yuan Rao, Lianwei Wu, and Ling Sun. 2022. Issues and challenges of aspect-based sentiment analysis: A comprehensive survey. *IEEE Transactions on Affective Computing*, 13(2):845–863.

Maria Pontiki, Dimitris Galanis, Haris Papageorgiou, Ion Androutsopoulos, Suresh Manandhar, Mohammad AL-Smadi, Mahmoud Al-Ayyoub, Yanyan Zhao, Bing Qin, Orphée De Clercq, Véronique Hoste, Marianna Apidianaki, Xavier Tannier, Natalia Loukachevitch, Evgeniy Kotelnikov, Nuria Bel, Salud María Jiménez-Zafra, and Gülşen Eryiğit. 2016. SemEval-2016 task 5: Aspect based sentiment analysis. In *Proceedings of the 10th International Workshop on Semantic Evaluation (SemEval-2016)*, pages 19–30, San Diego, California. Association for Computational Linguistics.

Maria Pontiki, Dimitris Galanis, Haris Papageorgiou, Suresh Manandhar, and Ion Androutsopoulos. 2015. SemEval-2015 task 12: Aspect based sentiment analysis. In *Proceedings of the 9th International Workshop on Semantic Evaluation (SemEval 2015)*, pages 486–495, Denver, Colorado. Association for Computational Linguistics.

Maria Pontiki, Dimitris Galanis, John Pavlopoulos, Harris Papageorgiou, Ion Androutsopoulos, and Suresh Manandhar. 2014. SemEval-2014 task 4: Aspect based sentiment analysis. In *Proceedings of the 8th International Workshop on Semantic Evaluation (SemEval 2014)*, pages 27–35, Dublin, Ireland. Association for Computational Linguistics.

Alec Radford, Jeffrey Wu, Rewon Child, David Luan, Dario Amodei, Ilya Sutskever, et al. 2019. Language models are unsupervised multitask learners. *OpenAI blog*, 1(8):9.

Seshan Ramaswami and Susheela Abraham Varghese. 2003. Reading the voice of the customer: A content analysis of consumer reviews.

Konstantinos I Roumeliotis, Nikolaos D Tselikas, and Dimitrios K Nasiopoulos. 2024. Llms in e-commerce: a comparative analysis of gpt and llama models in product review evaluation. *Natural Language Processing Journal*, 6:100056.

Marzieh Saeidi, Guillaume Bouchard, Maria Liakata, and Sebastian Riedel. 2016. SentiHood: Targeted aspect based sentiment analysis dataset for urban neighbourhoods. In *Proceedings of COLING 2016, the 26th International Conference on Computational Linguistics: Technical Papers*, pages 1546–1556, Osaka, Japan. The COLING 2016 Organizing Committee.

Hassan Saif, Miriam Fernández, Yulan He, and Harith Alani. 2013. Evaluation datasets for twitter sentiment analysis: a survey and a new dataset, the sts-gold. In *1st Interantional Workshop on Emotion and Sentiment in Social and Expressive Media: Approaches and Perspectives from AI (ESSEM 2013)*.

Taha Shangipour ataei, Kamyar Darvishi, Soroush Javdan, Behrouz Minaei-Bidgoli, and Sauleh Eetemadi. 2022. Pars-ABSA: a manually annotated aspect-based sentiment analysis benchmark on Farsi product reviews. In *Proceedings of the Thirteenth Language Resources and Evaluation Conference*, pages 7056–7060, Marseille, France. European Language Resources Association.

Frans Sudirjo, Karno Diantoro, Jassim Ahmad Al-Gasawneh, Hizbul Khootimah Azzaakiyyah, and Abu Muna Almaududi Ausat. 2023. Application of chatgpt in improving customer sentiment analysis for businesses. *Jurnal Teknologi Dan Sistem Informasi Bisnis*, 5(3):283–288.

Kian Long Tan, Chin Poo Lee, and Kian Ming Lim. 2023. A survey of sentiment analysis: Approaches, datasets, and future research. *Applied Sciences*, 13(7).

Rasika Wagh and Payal Punde. 2018. Survey on sentiment analysis using twitter dataset. In *2018 Second International Conference on Electronics, Communication and Aerospace Technology (ICECA)*, pages 208–211.

Mayur Wankhade, Annavarapu Chandra Sekhara Rao, and Chaitanya Kulkarni. 2022. A survey on sentiment analysis methods, applications, and challenges. *Artificial Intelligence Review*, 55(7):5731–5780.

Wenxuan Zhang, Yue Deng, Bing Liu, Sinno Pan, and Lidong Bing. 2024. Sentiment analysis in the era of large language models: A reality check. In *Findings of the Association for Computational Linguistics: NAACL 2024*, pages 3881–3906, Mexico City, Mexico. Association for Computational Linguistics.

Wenxuan Zhang, Xin Li, Yang Deng, Lidong Bing, and Wai Lam. 2022. A survey on aspect-based sentiment analysis: Tasks, methods, and challenges. *IEEE Trans. on Knowl. and Data Eng.*, 35(11):11019–11038.

## A   Prompts

### Listing 1: LLaMA-3 system prompt

```
messages = [ {"role": "system", "content
   ": """You are a helpful assistant
   that performs aspect based sentiment
    analysis for Hugo Boss! Do not
   communicate back, just provide the
   answer in the requested format"""},
   {"role": "user", "content": text}, ]
prompt = output.tokenizer.
   apply_chat_template( messages,
   tokenize=False,
   add_generation_prompt=True )
```

### Listing 2: GPT-4 system prompt

```
system_prompt = """You are a helpful
   assistant that performs aspect based
    sentiment analysis for Hugo Boss!
   Do not communicate back, just
   provide the answer in the requested
   format."""

prompt_value = StringPromptValue(text=
   chat_prompt_with_values)

output = llm.invoke([ SystemMessage(
   content=system_prompt), HumanMessage
   (content=prompt_value.text), ])
```

### Listing 3: LLaMA-3 user prompt

```
text = f"""First, identify the aspects
   in the provided review from the
   given list, and then find the
   customer sentiment (positive,
   neutral, or negative) for each of
   the aspect.
Make sure to take into account the
   difference in the language, cultural
    aspects, sarcasm, emojis, and other
    linguistic behaviours when
   interpreting and assessing the
   reviews.

Aspects: [Product (collection, item,
   quality, variety, display, selection
   ), Service (staff, assistance, crew,
    employee, attitude, handling,
   hospitality), Brand, Price, Store (
   shop, location, atmosphere), Online
   (order), Purchase, Return, General (
   overall shopping experience)]

You can follow the examples below:

Review: The product quality is great but
    the customer service is terrible.
Aspect: product
Sentiment: positive
Aspect: service
Sentiment: negative

Review: I love the location of the store
   . The collection and selection looks
    great, however the prices are too
   high.
```

```
Aspect: store
Sentiment: positive
Aspect: product
Sentiment: positive
Aspect: prices
Sentiment: negative

Review: All top, everything as I
   expected, recommend.
Aspect: general
Sentiment: positive

Now proceed with the following review:
   ```{review}```

Remember to strictly focus only on the
   aspects from the list and reply only
    with the answer in the following
   JSON format:
[{{"aspect1": "sentiment"}},
...
{{"aspectN": "sentiment"}}]
"""
```

### Listing 4: GPT-4 user prompt

```
First, identify the aspects in the
   provided review from the given list,
    and then find the customer
   sentiment (positive, neutral ot
   negative) for each of the aspect.

Aspects: [Product (collection, item,
   quality, variety, display), Service
   (staff, assistance, crew, employee,
   attitude, handling, hospitality),
   Brand, Price, Store (shop, location,
    atmosphere), Online (order), Return
   , General (overall shopping
   experience)]

You can follow the examples below:
review: "The␣product␣quality␣is␣great␣
   but␣the␣customer␣service␣is␣terrible
   ."
analysis: "[{{"product":␣"positive",␣"
   service":␣"negative"}}]"

review: "I␣love␣the␣location␣of␣the␣
   store.␣The␣collection␣looks␣great,␣
   however␣the␣prices␣are␣too␣high.",
analysis: "[{{"store":␣"positive",␣"
   product":␣"positive",␣"price":␣"
   negative"}}]"

review: "All␣top,␣everything␣as␣I␣
   expected,␣recommend.",
analysis: "[{{"general":␣"positive"}}]"

Make sure to take in account the
   difference in the language, cultural
    aspects, sarcasm, emojis and other
   linguistic behaviours when
   interpreting and assessing the
   reviews. Remember to strictly reply
   only with the answer in the
   following JSON format:
[{{"aspect1": "sentiment"}},
...
{{"aspectN": "sentiment"}}]
```

## B  Datasheet

### B.1  Motivation for Dataset Creation

**Why was the dataset created?** (e.g., were there specific tasks in mind, or a specific gap that needed to be filled?)

The dataset was created to enable aspect-based sentiment analysis on customer reviews using Large Language Models (LLMs).

**What (other) tasks could the dataset be used for?** Are there obvious tasks for which it should not be used?

The dataset could also be used for traditional sentiment analysis given it also contains star ratings for each review.

**Has the dataset been used for any tasks already?** If so, where are the results so others can compare (e.g., links to published papers)?

This paper is the first to use the dataset.

**Who funded the creation of the dataset?** If there is an associated grant, provide the grant number.

The data collection was supported by the consulting firm Metyis (`https://metyis.com/`).

### B.2  Dataset Composition

**What are the instances?** (that is, examples; e.g., documents, images, people, countries) Are there multiple types of instances? (e.g., movies, users, ratings; people, interactions between them; nodes, edges)

Each instance consists of a customer review for a brick-and-mortar store, scarped from Google maps.

**Are relationships between instances made explicit in the data (e.g., social network links, user-/movie ratings, etc.)?**

No.

**How many instances of each type are there?**

The dataset consists of 10,814 reviews.

**What data does each instance consist of?** "Raw" data (e.g., unprocessed text or images)? Features/attributes? Is there a label/target associated with instances? If the instances are related to people, are subpopulations identified (e.g., by age, gender, etc.) and what is their distribution?

In addition to the country and city of the store that is reviewed, the date the review was published at, its text, and the star rating are part of each instance. The language of the review is automatically annotated, while aspects and their sentiments have been manually annotated.

**Is everything included or does the data rely on external resources?** (e.g., websites, tweets, datasets) If external resources, a) are there guarantees that they will exist, and remain constant, over time; b) is there an official archival version. Are there licenses, fees or rights associated with any of the data?

Everything is included in the dataset.

**Are there recommended data splits or evaluation measures?** (e.g., training, development, testing; accuracy/AUC)

Since the dataset is designed for zero-shot classification, there is no recommended split. Given the imbalanced distribution of positive and negative sentiments, we recommend and evaluation measure that takes this into account, like F1-score.

**What experiments were initially run on this dataset?** Have a summary of those results and, if available, provide the link to a paper with more information here.

The dataset was initially used for the evaluation of the aspect-based sentiment analysis capabilities of LLMs.

### B.3  Data Collection Process

**How was the data collected?** (e.g., hardware apparatus/sensor, manual human curation, software program, software interface/API; how were these constructs/measures/methods validated?)

The data was collected using Apify (`https://apify.com/compass/google-maps-reviews-scraper`).

**Who was involved in the data collection process?** (e.g., students, crowdworkers) How were they compensated? (e.g., how much were crowdworkers paid?)

The data was collected by fully-qualified lawyers during their usual work-time. All participants worked for organizations that pay according to the collective labor agreement for public service workers in German states.

**Over what time-frame was the data collected?** Does the collection time-frame match the creation time-frame?

The data was collected in 2024. The reviews were written between 2012 and 2024.

**How was the data associated with each instance acquired?** Was the data directly observable (e.g., raw text, movie ratings), reported by subjects (e.g., survey responses), or indirectly inferred/derived from other data (e.g., part of speech tags; model-based guesses for age or language)? If

the latter two, were they validated/verified and if so how?

The reviews themselves and their metadata was directly observable, the language of the reviews was automatically derived using the Google Translate API, the aspects and their sentiments were manually annotated by the authors.

**Does the dataset contain all possible instances?** Or is it, for instance, a sample (not necessarily random) from a larger set of instances?

No, the dataset does not claim completeness in any sense.

**If the dataset is a sample, then what is the population?** What was the sampling strategy (e.g., deterministic, probabilistic with specific sampling probabilities)? Is the sample representative of the larger set (e.g., geographic coverage)? If not, why not (e.g., to cover a more diverse range of instances)? How does this affect possible uses?

The dataset spans multiple European countries and a time-frame of over a decade.

**Is there information missing from the dataset and why?** (this does not include intentionally dropped instances; it might include, e.g., redacted text, withheld documents) Is this data missing because it was unavailable?

Reviews that only consists of a star rating but do not provide any text have been excluded.

## B.4    Dataset Distribution

**How is the dataset distributed?** (e.g., website, API, etc.; does the data have a DOI; is it archived redundantly?)

It is archived on GitHub (`https://github.com/Responsible-NLP/ABSA-Retail-Corpus`).

**When will the dataset be released/first distributed?** (Is there a canonical paper/reference for this dataset?)

Publication of the paper.

**What license (if any) is it distributed under?** Are there any copyrights on the data?

The annotations are licensed under CC-BY-SA 4.0.

**Are there any fees or access/export restrictions?**

No.

## B.5    Dataset Maintenance

**Who is supporting/hosting/maintaining the dataset?** How does one contact the owner/curator/-manager of the dataset (e.g. email address, or other contact info)?

See the GitHub repository.

**Will the dataset be updated?** How often and by whom? How will updates/revisions be documented and communicated (e.g., mailing list, GitHub)? Is there an erratum?

There are no plans to update the dataset unless important mistakes become clear.

**If the dataset becomes obsolete how will this be communicated?**

On the GitHub page.

**Is there a repository to link to any/all papers/systems that use this dataset?**

Yes.

**If others want to extend/augment/build on this dataset, is there a mechanism for them to do so?** If so, is there a process for tracking/assessing the quality of those contributions. What is the process for communicating/distributing these contributions to users?

We would suggest to create a fork on GitHub.

## B.6    Legal & Ethical Considerations

**If the dataset relates to people (e.g., their attributes) or was generated by people, were they informed about the data collection?** (e.g., datasets that collect writing, photos, interactions, transactions, etc.)

There is no information about individuals in the data or was recorded during the annotation of the data.

**If it relates to other ethically protected subjects, have appropriate obligations been met?** (e.g., medical data might include information collected from animals)

N.a.

**If it relates to people, were there any ethical review applications/reviews/approvals?** (e.g. Institutional Review Board applications)

N.a.

**If it relates to people, were they told what the dataset would be used for and did they consent? What community norms exist for data collected from human communications?** If consent was obtained, how? Were the people provided with any mechanism to revoke their consent in the future or for certain uses?

N.a.

**If it relates to people, could this dataset expose people to harm or legal action?** (e.g., financial social or otherwise) What was done to mitigate or reduce the potential for harm?

N.a.

**If it relates to people, does it unfairly advantage or disadvantage a particular social group?** In what ways? How was this mitigated?

N.a.

**If it relates to people, were they provided with privacy guarantees?** If so, what guarantees and how are these ensured?

N.a.

**Does the dataset comply with the EU General Data Protection Regulation (GDPR)?** Does it comply with any other standards, such as the US Equal Employment Opportunity Act?

Yes, since only publicly available information was collected, the dataset complies with the GDPR and similar regulations.

**Does the dataset contain information that might be considered sensitive or confidential?** (e.g., personally identifying information)

No.

**Does the dataset contain information that might be considered inappropriate or offensive?**

No.

# The Need for Robust and Inclusive Benchmarks in Evaluating LLMs on Arabic Text

**Lubana Al Rayes**
Department of Computer Science
University of Sharjah
Sharjah, UAE
lrayes@sharjah.ac.ae

**Ashraf Elnagar**
Department of Computer Science
University of Sharjah
Sharjah, UAE
ashraf@sharjah.ac.ae

## Abstract

The widespread success of large language models (LLMs) has prompted increasing interest in their evaluation across diverse linguistic settings, yet systematic assessments for Arabic remain underexplored. This survey presents a structured taxonomy of benchmarks specifically designed to evaluate LLMs on Arabic text. It critically reviews existing benchmarks, highlighting their coverage across multiple domains, including general single-task and multi-task scenarios, knowledge and reasoning tasks, and domain-specific applications. Finally, it identifies key methodological limitations and proposes future research directions to facilitate the development of more robust, inclusive, and culturally aligned evaluation frameworks for LLMs.

## 1 Introduction

Large Language Models (LLMs) have become a cornerstone of modern natural language processing (NLP), demonstrating remarkable performance across a wide spectrum of tasks such as machine translation (MT), sentiment analysis, dialogue generation, and reasoning (Yang et al., 2025). Their broad generalization capabilities have positioned them as foundational tools in diverse domains, ranging from healthcare and law to education and creative writing (Bommasani et al., 2021). However, their widespread deployment necessitates rigorous evaluation frameworks to ensure reliability, fairness, and robust performance in complex reasoning, factual consistency, and linguistic competence, particularly in low-resource languages like Arabic.

Although Arabic is among the most widely spoken languages globally, it is significantly underrepresented in the training data of many multilingual large language models (MLLMs), where English typically accounts for over 90% of the corpus and Arabic often constitutes less than 1% (Xu et al., 2025; Qian et al., 2024). Consequently,

many Arabic-centric or multilingual models struggle to maintain consistent performance across dialects, linguistic styles, and culturally grounded tasks (Magdy et al., 2025; Alwajih et al., 2025). To address this gap, an increasing number of benchmarks have been proposed to evaluate LLMs on Arabic tasks. These benchmarks span a variety of domains and evaluation objectives, including general multi-task performance, commonsense and factual reasoning, domain-specific applications (e.g., legal and healthcare), and fine-grained single-task assessments. Despite this growing body of work, there is no unified or comprehensive framework that consolidates these efforts to guide comparative evaluation or diagnostic analysis.

This paper addresses these challenges by offering a structured survey focused exclusively on benchmarks used to evaluate LLMs on Arabic text. It systematically reviews existing benchmarks and organizes them into a unified taxonomy based on task type and domain focus. The paper also identifies common methodological gaps and proposes directions for future research. This survey serves as a foundational resource for researchers and practitioners seeking to understand the current landscape, design more inclusive benchmarks, or select appropriate evaluation frameworks for their models.

## 2 Related Work

Several recent surveys have synthesized progress in LLM development and evaluation, yet none have specifically focused on existing benchmarks for evaluating LLMs on Arabic text.

One of the most relevant works is by Mashaabi et al. (2025). This survey provides an overview of Arabic LLMs across different architectures (encoder-only, decoder-only, encoder-decoder), linguistic forms (Modern Standard Arabic (MSA), Classical Arabic, Dialectal Arabic), and pretraining datasets. It also evaluates the openness of these models and their performance across downstream

NLP tasks. However, the work does not systematically survey evaluation benchmarks used to assess these models. Benchmarks are only briefly mentioned in the context of task-based performance. In related efforts, benchmarks focusing specifically on Arabic word embeddings and contextualized embeddings have been proposed, including those by Yagi et al. (2023) and Elnagar et al. (2023), providing comprehensive evaluation frameworks for these foundational models. Furthermore, studies examining Arabic punctuation and its linguistic characteristics have offered insights into its rule-governed nature (Yagi et al., 2024).

Similarly, Rhel and Roussinov (2025) offer a general overview of Arabic LLMs. While the paper reflects on limitations in Arabic resources and the application of LLMs to Arabic NLP tasks, its focus is not on benchmarking. Instead, it summarizes the adoption of LLMs in Arabic contexts and briefly lists common datasets, without detailed analysis or categorization of benchmarks used across tasks. On a related note, cross-lingual models integrating Arabic language with images have recently been developed, such as the AraCLIP framework by Al-Barham et al. (2025), which explores novel approaches to Arabic vision-language understanding.

Outside the Arabic context, Laskar et al. (2024) presented a systematic review of LLM evaluation pipelines, identifying challenges such as reproducibility, dataset contamination, and fairness across benchmarks. Their work offers a robust foundation for understanding the complexities of LLM evaluation but focuses primarily on English and multilingual settings. Likewise, Lai et al. (2023) analyzed the multilingual performance of ChatGPT across 37 languages, including Arabic, through zero-shot evaluations on tasks like summarization and Part of Speech (POS) tagging. While their work evaluated Arabic among other languages, it did not aim to survey benchmarks nor did it focus on Arabic text.

To the best of our knowledge, this paper is the first to focus specifically on the evaluation benchmarks used to assess LLMs on Arabic text, rather than surveying Arabic LLMs themselves. While prior surveys have examined Arabic language models in terms of architecture, datasets, and application domains, none have systematically analyzed the benchmarks that underpin their evaluation. This distinction allows our work to fill a critical gap by offering a structured overview of the evaluation landscape and identifying methodological shortcomings in current benchmarking practices.

## 3 Methodology

A total of 26 relevant studies were included in this survey paper. All studies were published between 2022 and 2025. The search window spanned 2020 to 2025, and the methodology followed a systematic approach structured into three main phases.

### 3.1 Literature search

To identify relevant research, a comprehensive literature search was conducted across multiple scientific databases, including Google Scholar, Elsevier, and IEEE Xplore. The search queries used combinations of keywords such as "Large Language Models", "Benchmark", "Evaluation", and "Arabic Text". This process yielded a total of 42 records.

### 3.2 Inclusion and exclusion criteria

The retrieved records were screened for eligibility using predefined criteria. Studies were included if they evaluated LLMs on Arabic text, regardless of whether other languages were involved, provided Arabic evaluation was a core component. Studies that focused exclusively on non-textual modalities (e.g., images, audio, video) or did not contain Arabic content were excluded.

Duplicates were identified and removed (n = 2), resulting in 40 records screened. Of these, 4 were excluded during initial screening due to irrelevance, and 10 more were excluded after full-text assessment. No records were missing or unretrievable. In total, 26 unique studies met the inclusion criteria and were included in the final review. A visual summary of this process is shown in Figure 1.

### 3.3 Taxonomy

The selected studies were organized using a structured taxonomy designed to categorize LLM evaluations on Arabic text. Each study was assigned to a distinct subcategory under one of four main categories, based on its primary objective and evaluation scope. The taxonomy comprises:

- General Multi-Task Evaluation Benchmarks

- Knowledge and Reasoning Benchmarks

- Domain-Specific Benchmarks

- Focused Single-Task Evaluations

Figure 1: Flow-chart for study inclusion

Subcategories reflect the evaluation scope, task specificity, and domain orientation of each study, as illustrated in Figure 2 and detailed in the following sections.

While our taxonomy was initially designed around the specific context of LLM evaluations on Arabic text, its fundamental structure is language-agnostic and can be generalized across diverse linguistic contexts, potentially serving as a broader blueprint for evaluating LLMs.

## 4 Taxonomy for Evaluating LLMs on Arabic Text

This section presents the taxonomy used to classify benchmarks for evaluating LLMs on Arabic text. The taxonomy is divided into four major categories: (1) General Multi-Task Evaluation Benchmarks, (2) Knowledge and Reasoning Benchmarks, (3) Domain-Specific Benchmarks, and (4) Focused Single-Task Evaluations. Each category captures distinct evaluation objectives, methodological designs, and linguistic considerations.

Detailed characteristics of each benchmark are summarized in Appendix A.

### 4.1 General Multi-Task Evaluation Benchmarks

General multi-task evaluation benchmarks are designed to assess LLMs on a broad range of NLP tasks that combine natural language understanding (NLU) and generation (NLG). Within this category, we distinguish between two subcategories: Multi-Task Mixed NLU/NLG Benchmarks and NLG-

Focused Multi-Task Benchmarks.

The first subcategory, Multi-Task Mixed NLU/NLG Benchmarks, includes benchmarks that evaluate LLMs across diverse general-domain tasks. One example is the AraT5/ARGEN benchmark (Elmadany et al., 2022), which adopts a text-to-text format to uniformly structure input and output for eight tasks, including sentiment analysis, classification, Named Entity Recognition (NER), extractive QA, summarization, and paraphrasing. The benchmark tests models like AraT5, mT5, and mBART in zero- and few-shot settings, using task-appropriate metrics such as F1, BLEU, and ROUGE. Despite its extensive task coverage, the benchmark is primarily based on MSA, with minimal attention to dialectal Arabic. This limits its applicability in real-world scenarios involving linguistic variation.

Another benchmark in this subcategory is GPTAraEval (Khondaker et al., 2023), which assesses ChatGPT-3.5 and GPT-4 across 44 tasks drawn from 60 datasets, encompassing classification, paraphrase detection, QA, and NER. The benchmark operates exclusively in zero-shot mode to reflect typical usage of proprietary LLMs. While GPT-4 demonstrates superior performance over its predecessor, the study's focus on only two models introduces bias and excludes insights from Arabic-centric or fine-tuned models.

LAraBench (Abdelali et al., 2023) expands multi-task evaluation by including speech-related tasks, such as automatic speech recognition (ASR) and text-to-speech (TTS), in addition to standard NLP tasks. It covers 33 tasks across 61 datasets and evaluates models including GPT-4, Jais, and Whisper. The benchmark shows that even the strongest LLMs face difficulties with syntactic and sequence tagging tasks. These issues are partly due to the lack of Arabic-specific pretraining and inconsistent output formatting. The models also perform poorly across different Arabic language varieties, which can be attributed to the lack of dialectal data.

The second subcategory, NLG-Focused Multi-Task Benchmarks, centers specifically on generative language capabilities. The Dolphin benchmark (Elmadany et al., 2023) exemplifies this by focusing exclusively on Arabic NLG tasks, including summarization, storytelling, dialogue, and data-to-text generation. Comprising 200,000 completions across 20,000 prompts, Dolphin evaluates LLMs like GPT-4, Falcon, and ChatGPT using both hu-

Figure 2: Taxonomy for Evaluating LLMs on Arabic Text

man judgments (e.g., grammaticality, coherence) and automatic metrics (e.g., BLEU, ROUGE-L, COMET). While Dolphin provides a rich resource for assessing generative fluency and factuality, a limitation of considering only NLG is that it overlooks other critical language understanding capabilities, such as reasoning, retrieval, and classification.

Benchmarks under the general multi-task category offer foundational insights into the capabilities of LLMs in Arabic across diverse tasks. However, limitations such as restricted dialectal coverage, model scope, and narrow task focus indicate a need for more comprehensive, balanced, and culturally representative evaluation frameworks.

## 4.2 Knowledge and Reasoning Benchmarks

Knowledge and reasoning benchmarks aim to assess the depth of logical inference and factual understanding of LLMs beyond basic comprehension. These are typically structured as multi-choice questions (MCQs) or explanatory tasks designed to simulate complex, real-world problem-solving situations.

A primary subcategory is Massive Multitask QA Benchmarks, which assess a model's breadth of knowledge across subjects. For example, ArabicMMLU (Koto et al., 2024) covers 14,575 MCQs across 40 tasks, drawing from real-world school exams in various Arabic-speaking regions. Similarly, AlGhafa (Almazrouei et al., 2023) includes 7,226 MCQs across 45 tasks, categorized into reasoning, knowledge, reading comprehension, and math. Another example is AraSTEM (Mustapha et al., 2024),

which focuses on STEM subjects with over 11,000 questions ranging from primary school to college-level. Finally, the Qiyas Benchmark (Al-Khalifa and Al-Khalifa, 2024) evaluates models using questions from the Saudi General Aptitude Test, covering both verbal and mathematical reasoning. These benchmarks offer broad task coverage, but their formats rely entirely on MCQs, which simplify the task structure and may inflate performance by enabling guessing (Koto et al., 2024; Almazrouei et al., 2023; Mustapha et al., 2024; Al-Khalifa and Al-Khalifa, 2024). Such format constraints can limit a model's opportunity to demonstrate deeper reasoning or generative capabilities. Additionally, evaluating only a narrow set of models restricts the ability to offer a comprehensive view of performance across the broader LLM landscape (Al-Khalifa and Al-Khalifa, 2024), including emerging or open-source models. Most benchmarks are also confined to MSA, excluding dialects, informal text, or culturally specific content (Koto et al., 2024; Mustapha et al., 2024; Al-Khalifa and Al-Khalifa, 2024).

The second subcategory, Commonsense Reasoning Benchmarks, evaluates a model's intuitive understanding of everyday scenarios. ArabicSense (Lamsiyah et al., 2025) is a newly proposed benchmark that assesses commonsense validation, explanation selection, and generative explanation. The dataset is synthetically generated and covers a range of reasoning skills. However, it remains limited in scope, focusing only on three task types and lacking the diversity of real-world language use.

Additionally, its synthetic nature may introduce biases or overfitting tendencies not representative of actual human-authored content.

### 4.3 Domain-Specific Benchmarks

Domain-specific benchmarks are designed to evaluate LLMs on tasks rooted in real-world applications and specialized knowledge areas. Unlike general-purpose benchmarks, which assess broad linguistic competence, these benchmarks target specific domains, such as law, health, cultural reasoning, and safety, to assess how well models handle context-sensitive, factual, and domain-relevant language use. This subsection is organized into four sub-categories of domain-specific benchmarks: legal, cultural and dialectal competence, health, and trustworthiness and safety.

In the legal domain, the ArabLegalEval benchmark (Hijazi et al., 2024) provides a multi-task framework designed to evaluate Arabic LLMs' legal reasoning capabilities. It includes over 15,000 instances covering MCQs, open-ended QA, and carefully translated items from the English-language LegalBench dataset. These tasks primarily draw from Saudi legal sources, such as regulations on consumer contracts and privacy policies. While ArabLegalEval provides a rigorous and diverse evaluation setting, its heavy reliance on Saudi legal texts may limit its applicability across broader Arabic legal systems.

Cultural and dialectal competence has emerged as a critical dimension in evaluating LLMs on Arabic text due to the region's linguistic diversity. AraDiCE (Mousi et al., 2025) benchmarks dialectal and cultural understanding across Egyptian, Gulf, Levantine, and MSA. It spans dialect identification, misinformation detection, and cultural reasoning. However, it primarily relies on synthetic data generated via machine translation with post-editing, which may introduce unnatural phrasing.In addition, the omission of key dialects such as Maghrebi limits its regional coverage. The Palm benchmark (Alwajih et al., 2025) offers 17,411 annotated instruction–response pairs covering ten dialects across 20 culturally salient domains. Despite its breadth, Palm exhibits skewed country-level representation. Similarly, the SaudiCulture benchmark (Ayash et al., 2025) evaluates LLMs on region-specific cultural questions within Saudi Arabia, capturing intranational differences across five regions. Nonetheless, its geographic scope limits

generalizability to broader Arab cultural contexts. Jawaher benchmark (Magdy et al., 2025) targets proverb translation and explanation in 20 dialects, exposing the limitations of current LLMs in handling idiomatic, figurative, and culturally grounded expressions. However, its evaluation is affected by the use of English-only prompts, which limits the assessment of models' native Arabic comprehension. Lastly, the culturally aligned benchmark (Nacar et al., 2025) critiques the Western bias of traditional evaluation frameworks and introduces ILMAAM, a curated leaderboard tailored to Arabic sociocultural contexts. It improves cultural appropriateness.

In the health domain, the Health Claims benchmark (obaid Alharbi et al., 2025) evaluates GPT-4's ability to classify and verify health-related claims across Saudi, Egyptian, Lebanese, and Moroccan dialects. The study utilizes 329 expert-verified claims from AraFacts and ArCOV19-Rumors, generating 6,520 dialect-specific queries with varying presupposition levels. It applies a novel Cultural Sensitivity Score to measure context-aware accuracy. The benchmark is limited by its evaluation of only a single model (GPT-4), which restricts its comparative utility, and by its narrow dialectal coverage that excludes other widely spoken Arabic varieties.

The domain of trustworthiness and safety is addressed by AraTrust (Alghamdi et al., 2025), which includes 522 multiple-choice questions evaluating LLMs on ethics, legality, offensiveness, and privacy. It introduces evaluations across several prompting settings, including chain-of-thought reasoning. However, the benchmark's exclusive use of multiple-choice formats restricts deeper assessment of models' ethical reasoning in open-ended contexts.

### 4.4 Focused Single-Task Evaluations

Benchmarks in this category are designed to evaluate LLMs on narrowly defined tasks that test specific competencies in Arabic. Unlike multi-task benchmarks, these evaluations isolate a single task, such as sentiment classification, machine translation, or hallucination detection, allowing for more fine-grained assessment of model performance. This category comprises three major sub-categories of tasks: classification and understanding, generation and transformation, and factuality and hallucination detection.

Classification and Understanding Tasks target the ability of LLMs to label and disambiguate text based on semantic, syntactic, or pragmatic cues. In the domain of sentiment classification, Al-Thubaity et al. (2023) evaluated GPT-3.5, GPT-4, and PaLM 2 (Bard AI) using the Saudi Dialect Twitter Corpus, covering a small-scale dataset of 2,690 tweets labeled as positive, negative, or neutral. The benchmark revealed close performance between GPT-4 and fine-tuned BERT baselines, yet it is restricted to a single dialect (Saudi). A benchmark for Cross-Lingual NER was proposed by Al-Duwais et al. (2024) to test six multilingual LLMs using seven datasets across domains like news and social media. The benchmark revealed strong performance by encoder-based models such as XLM-R and mBERT. Abdel-Salam (2024) introduced a benchmark for Word Sense and Location Mention Disambiguation using SALMA and IDRISI-D datasets. While demonstrating LLM competence in controlled zero-shot setups, the benchmark omits dialectal variations and depends heavily on short contexts and English translations for retrieval.

Generation and Transformation Tasks evaluate how well models perform structured text transformation, such as translation or correction. A machine translation benchmark was proposed by Kadaoui et al. (2023) using 1,000 dialectal Arabic sentences across ten varieties from the MADAR corpus. While the benchmark spans several dialects, it includes only two LLMs, ChatGPT and Bard. Kwon et al. (2023) benchmarked LLMs on Arabic grammatical error correction (AGEC) using QALB datasets, evaluating performance with prompting strategies such as zero-shot, few-shot, and instruction tuning. The benchmark highlights LLM underperformance on semantic errors and lacks dialectal diversity. Another example in this group is the punctuation restoration benchmark (Al Wazrah et al., 2025), which uses a curated dataset of 10,046 paragraphs to test seven LLMs and a fine-tuned AraBERT. GPT4-o performed best overall. The benchmark suffers from skewed punctuation distributions.

Factuality and hallucination evaluation tasks assess LLMs' ability to distinguish between true and false claims or to avoid generating fabricated content. Gupta et al. (2025) developed a fact-checking benchmark using 771 claims from the X-Fact dataset, focusing on binary classification with English reasoning strategies applied to Arabic

input. The dataset is heavily skewed toward false claims and excludes recent advanced models like GPT-4, which limits longitudinal comparisons. In the area of hallucination detection, the Halwasa benchmark (Mubarak et al., 2024) evaluates Arabic hallucinations using 10,000 synthetic factual sentences generated by LLMs. The dataset was created using 1,000 randomly selected words from the SAMER Arabic readability lexicon. For each word, both GPT-3.5 and GPT-4 were prompted to generate ten factual Arabic sentences. After filtering out duplicates and invalid outputs, five unique sentences were retained per model, resulting in 5,000 sentences from each and a total of 10,000 sentences. Each sentence was manually annotated by trained human annotators across four dimensions: (1) whether it makes a verifiable factual claim, (2) whether the claim is factually correct, (3) whether the sentence follows proper Arabic grammar, and (4) the reference sources used for factual verification. A key limitation of this benchmark is its exclusive focus on just two models, GPT-3.5 and GPT-4, which restricts its comparative scope across a broader range of Arabic or multilingual LLMs. Similarly, HalluVerse25 (Abdaljalil et al., 2025) is a multilingual hallucination detection benchmark that includes 828 Arabic sentence pairs focused on biographical content. While it supports cross-lingual comparison, the benchmark inherits potential biases from Wikidata and the use of GPT-generated data, constraining its generalizability beyond the biographical domain.

## 5 Critical Analysis of Existing Arabic LLM Evaluations

Despite significant advancements in evaluating LLMs for Arabic text, existing benchmarks reveal several critical challenges. Specifically, there is a pronounced absence of separate intrinsic and extrinsic evaluations. Currently, benchmarks frequently blend these tasks into general multi-task evaluations, making it difficult to comprehensively assess specific competencies such as linguistic understanding, factual reasoning, and cultural awareness. This methodological conflation fails to provide a clear diagnostic of a model's performance, particularly in distinguishing whether success is driven by deep comprehension or surface-level task handling.

Another considerable limitation lies in the limited scope of model evaluations. Most benchmarks evaluate only a narrow set of LLMs, predominantly

focusing on well-known models such as GPT variants, neglecting emerging or specialized Arabic-centric models. Consequently, this narrow selection restricts the ability to address crucial comparative questions, such as identifying which models excel in specific tasks. Moreover, the scarcity of comparative analyses across a broader spectrum of models limits insights into model scalability and adaptability in diverse Arabic linguistic environments.

Additionally, the prevalent reliance exclusively on MCQs in many benchmarks represents another critical limitation. Solely using MCQs inherently simplifies evaluation tasks, potentially inflating model performance by allowing for guessing and limiting the ability to assess more sophisticated generative or explanatory capabilities.

In parallel with these methodological considerations, it is equally important to situate Arabic within the broader multilingual evaluation landscape. While this survey focuses on Arabic benchmarks for evaluating LLMs, understanding how Arabic is represented across cross-lingual benchmarks provides valuable context. Several cross-lingual benchmarks, such as XTREME (Hu et al., 2020), XGLUE (Liang et al., 2020), Blend (Myung et al., 2024), and (Chollampatt et al., 2025), include Arabic alongside other languages, often as a representative of Semitic or low-resource linguistic groups. However, these benchmarks typically offer limited task coverage for Arabic and rarely account for the linguistic diversity within the language, such as dialectal variation or cultural specificity. In contrast, Arabic-specific benchmarks provide more fine-grained evaluations tailored to the complexities of Arabic, including dialect identification, cultural reasoning, and script variants. Moreover, while cross-lingual benchmarks are valuable for assessing generalization and transfer learning, they often rely on translated or parallel data that may not reflect authentic language use. Arabic-centric benchmarks, by contrast, frequently involve native-authored content and culturally grounded tasks, offering a more accurate assessment of LLM performance on Arabic.

In addition to broadening evaluation contexts, this survey primarily focuses on benchmarking coverage and evaluation frameworks, we acknowledge the importance of analyzing bias in LLMs more explicitly. Several Arabic benchmarks, such as Ara-Trust and Palm, begin to address dimensions of bias

related to ethics, offensiveness, and regional representation. However, most existing datasets lack systematic annotations for sensitive attributes like gender, dialect, or sociopolitical context, making it difficult to assess fairness across subpopulations. Furthermore, benchmarks that rely on machine-translated or synthetic data may introduce unintended cultural or linguistic biases.

## 6 Conclusion and Future Directions

This survey has provided a comprehensive overview of existing benchmarks for evaluating LLMs on Arabic text, highlighting both significant progress and critical gaps. While current benchmarks offer valuable insights across various linguistic tasks and domains, they often conflate intrinsic and extrinsic evaluations, focus narrowly on a limited set of popular models, and rely heavily on simplified formats such as multiple-choice questions. Moreover, the underrepresentation of Arabic dialects and cultural nuances limits the applicability of these evaluations to the diverse Arabic language landscape. Bias and fairness considerations remain insufficiently addressed in most datasets, posing challenges for equitable model assessment.

To advance the field, future research should explicitly differentiate intrinsic language-specific evaluations (e.g., syntactic parsing, semantic understanding, morphological analysis) from extrinsic task-based assessments focused on real-world applications such as healthcare, law, and education. Expanding the range of evaluated models to include emerging, open-source, and Arabic-centric LLMs will enhance comparative analyses and foster innovation tailored to Arabic's unique linguistic characteristics.

Future benchmarks must incorporate diverse, realistic datasets reflecting dialectal variety and cultural context to improve real-world relevance. The growing importance of prompt engineering calls for systematic exploration of prompt formulations in both Arabic and English to optimize model performance and reliability. Additionally, incorporating bias-sensitive design principles and targeted fairness metrics is essential to ensure equitable evaluation across dialects, regions, and sociolinguistic groups.

Overall, addressing these methodological and practical gaps will deepen understanding of how LLMs perform on Arabic text and guide the development of more robust, culturally aware, and

effective language technologies.

## Limitations

This survey is limited by its exclusive focus on publicly documented academic benchmarks, omitting proprietary or industrial evaluations that may provide additional perspectives.

## References

Samir Abdaljalil, Hasan Kurban, and Erchin Serpedin. 2025. Halluverse25: Fine-grained multilingual benchmark dataset for llm hallucinations.

Reem Abdel-Salam. 2024. rematchka at arabicnlu2024: Evaluating large language models for arabic word sense and location sense disambiguation. In *Proceedings of the Second Arabic Natural Language Processing Conference*, pages 383–392.

Ahmed Abdelali, Hamdy Mubarak, Shammur Absar Chowdhury, Maram Hasanain, Basel Mousi, Sabri Boughorbel, Yassine El Kheir, Daniel Izham, Fahim Dalvi, Majd Hawasly, et al. 2023. Larabench: Benchmarking arabic ai with large language models. *arXiv preprint arXiv:2305.14982*.

Muhammad Al-Barham, Imad Afyouni, Khalid Almubarak, Ayad Turky, Ibrahim Abaker Targio Hashem, Ali Bou Nassif, Ismail Shahin, and Ashraf Elnagar. 2025. Unlocking language boundaries: Araclip-transforming arabic language and image understanding through cross-lingual models. *Engineering Applications of Artificial Intelligence*, 151:110577.

Mashael Al-Duwais, Hend Al-Khalifa, and Abdulmalik Al-Salman. 2024. A benchmark evaluation of multilingual large language models for arabic cross-lingual named-entity recognition. *Electronics*, 13(17):3574.

Shahad Al-Khalifa and Hend Al-Khalifa. 2024. The qiyas benchmark: Measuring chatgpt mathematical and language understanding in arabic. In *Proceedings of the 7th International Conference on Natural Language and Speech Processing (ICNLSP 2024)*, pages 343–351.

Abdulmohsen Al-Thubaity, Sakhar Alkhereyf, Hanan Murayshid, Nouf Alshalawi, Maha Omirah, Raghad Alateeq, Rawabi Almutairi, Razan Alsuwailem, Manal Alhassoun, and Imaan Alkhanen. 2023. Evaluating chatgpt and bard ai on arabic sentiment analysis. In *Proceedings of ArabicNLP 2023*, pages 335–349.

Asma Ali Al Wazrah, Afrah Altamimi, Hawra Aljasim, Waad Alshammari, Rawan Al-Matham, Omar Elnashar, Mohamed Amin, and Abdulrahman AlOsaimy. 2025. Evaluation of large language models on arabic punctuation prediction. In *Proceedings of the 1st Workshop on NLP for Languages Using Arabic Script*, pages 144–154.

Emad A Alghamdi, Reem Masoud, Deema Alnuhait, Afnan Y Alomairi, Ahmed Ashraf, and Mohamed Zaytoon. 2025. Aratrust: An evaluation of trustworthiness for llms in arabic. In *Proceedings of the 31st International Conference on Computational Linguistics*, pages 8664–8679.

Ebtesam Almazrouei, Ruxandra Cojocaru, Michele Baldo, Quentin Malartic, Hamza Alobeidli, Daniele Mazzotta, Guilherme Penedo, Giulia Campesan, Mugariya Farooq, Maitha Alhammadi, et al. 2023. Alghafa evaluation benchmark for arabic language models. In *Proceedings of ArabicNLP 2023*, pages 244–275.

Fakhraddin Alwajih, Abdellah El Mekki, Samar Mohamed Magdy, Abdelrahim A Elmadany, Omer Nacar, El Moatez Billah Nagoudi, Reem Abdel-Salam, Hanin Atwany, Youssef Nafea, Abdulfattah Mohammed Yahya, et al. 2025. Palm: A culturally inclusive and linguistically diverse dataset for arabic llms. *arXiv preprint arXiv:2503.00151*.

Lama Ayash, Hassan Alhuzali, Ashwag Alasmari, and Sultan Aloufi. 2025. Saudiculture: A benchmark for evaluating large language models cultural competence within saudi arabia. *arXiv preprint arXiv:2503.17485*.

Rishi Bommasani, Drew A Hudson, Ehsan Adeli, Russ Altman, Simran Arora, Sydney von Arx, Michael S Bernstein, Jeannette Bohg, Antoine Bosselut, Emma Brunskill, et al. 2021. On the opportunities and risks of foundation models. *arXiv preprint arXiv:2108.07258*.

Shamil Chollampatt, Minh Quang Pham, Sathish Reddy Indurthi, and Marco Turchi. 2025. Cross-lingual evaluation of multilingual text generation. In *Proceedings of the 31st International Conference on Computational Linguistics*, pages 7766–7777.

AbdelRahim Elmadany, Muhammad Abdul-Mageed, et al. 2022. Arat5: Text-to-text transformers for arabic language generation. In *Proceedings of the 60th annual meeting of the association for computational linguistics (Volume 1: Long papers)*, pages 628–647.

Abdelrahim Elmadany, Ahmed El-Shangiti, Muhammad Abdul-Mageed, et al. 2023. Dolphin: A challenging and diverse benchmark for arabic nlg. In *Findings of the Association for Computational Linguistics: EMNLP 2023*, pages 1404–1422.

Ashraf Elnagar, Sane Yagi, Youssef Mansour, Leena Lulu, and Shehdeh Fareh. 2023. A benchmark for evaluating arabic contextualized word embedding models. *Information Processing & Management*, 60(5):103452.

Ayushman Gupta, Aryan Singhal, Thomas Law, Veekshith Rao, Evan Duan, and Ryan Luo Li. 2025. Can llms verify arabic claims? evaluating the arabic fact-checking abilities of multilingual llms. In *Proceedings of the 1st Workshop on NLP for Languages Using Arabic Script*, pages 104–113.

Faris Hijazi, Somayah Alharbi, Abdulaziz AlHussein, Harethah Shairah, Reem Alzahrani, Hebah Alshamlan, George Turkiyyah, and Omar Knio. 2024. Arablegaleval: A multitask benchmark for assessing arabic legal knowledge in large language models. In *Proceedings of The Second Arabic Natural Language Processing Conference*, pages 225–249.

Junjie Hu, Sebastian Ruder, Aditya Siddhant, Graham Neubig, Orhan Firat, and Melvin Johnson. 2020. Xtreme: A massively multilingual multi-task benchmark for evaluating cross-lingual generalisation. In *International conference on machine learning*, pages 4411–4421. PMLR.

Karima Kadaoui, Samar M Magdy, Abdul Waheed, Md Tawkat Islam Khondaker, Ahmed Oumar El-Shangiti, El Moatez Billah Nagoudi, and Muhammad Abdul-Mageed. 2023. Tarjamat: Evaluation of bard and chatgpt on machine translation of ten arabic varieties. *arXiv preprint arXiv:2308.03051*.

Md Tawkat Islam Khondaker, Abdul Waheed, El Moatez Billah Nagoudi, and Muhammad Abdul-Mageed. 2023. Gptaraeval: A comprehensive evaluation of chatgpt on arabic nlp. *arXiv preprint arXiv:2305.14976*.

Fajri Koto, Haonan Li, Sara Shatnawi, Jad Doughman, Abdelrahman Boda Sadallah, Aisha Alraeesi, Khalid Almubarak, Zaid Alyafeai, Neha Sengupta, Shady Shehata, et al. 2024. Arabicmmlu: Assessing massive multitask language understanding in arabic. *arXiv preprint arXiv:2402.12840*.

Sang Yun Kwon, Gagan Bhatia, El Moatez Billah Nagoudi, and Muhammad Abdul-Mageed. 2023. Beyond english: Evaluating llms for arabic grammatical error correction. *arXiv preprint arXiv:2312.08400*.

Viet Lai, Nghia Ngo, Amir Pouran Ben Veyseh, Hieu Man, Franck Dernoncourt, Trung Bui, and Thien Nguyen. 2023. Chatgpt beyond english: Towards a comprehensive evaluation of large language models in multilingual learning. In *Findings of the Association for Computational Linguistics: EMNLP 2023*, pages 13171–13189.

Salima Lamsiyah, Kamyar Zeinalipour, Matthias Brust, Marco Maggini, Pascal Bouvry, Christoph Schommer, et al. 2025. Arabicsense: A benchmark for evaluating commonsense reasoning in arabic with large language models. In *Proceedings of the 4th Workshop on Arabic Corpus Linguistics (WACL-4)*, pages 1–11.

Md Tahmid Rahman Laskar, Sawsan Alqahtani, M Saiful Bari, Mizanur Rahman, Mohammad Abdullah Matin Khan, Haidar Khan, Israt Jahan, Amran Bhuiyan, Chee Wei Tan, Md Rizwan Parvez, et al. 2024. A systematic survey and critical review on evaluating large language models: Challenges, limitations, and recommendations. In *Proceedings of the 2024 Conference on Empirical Methods in Natural Language Processing*, pages 13785–13816.

Yaobo Liang, Nan Duan, Yeyun Gong, Ning Wu, Fenfei Guo, Weizhen Qi, Ming Gong, Linjun Shou, Daxin Jiang, Guihong Cao, et al. 2020. Xglue: A new benchmark dataset for cross-lingual pre-training, understanding and generation. *arXiv preprint arXiv:2004.01401*.

Samar M Magdy, Sang Yun Kwon, Fakhraddin Alwajih, Safaa Abdelfadil, Shady Shehata, and Muhammad Abdul-Mageed. 2025. Jawaher: A multidialectal dataset of arabic proverbs for llm benchmarking. *arXiv preprint arXiv:2503.00231*.

Malak Mashaabi, Shahad Al-Khalifa, and Hend Al-Khalifa. 2025. A survey of large language models for arabic language and its dialects.

Basel Mousi, Nadir Durrani, Fatema Ahmad, Md. Arid Hasan, Maram Hasanain, Tameem Kabbani, Fahim Dalvi, Shammur Absar Chowdhury, and Firoj Alam. 2025. AraDiCE: Benchmarks for dialectal and cultural capabilities in LLMs. In *Proceedings of the 31st International Conference on Computational Linguistics*, pages 4186–4218, Abu Dhabi, UAE. Association for Computational Linguistics.

Hamdy Mubarak, Hend Al-Khalifa, and Khaloud Suliman Alkhalefah. 2024. Halwasa: Quantify and analyze hallucinations in large language models: Arabic as a case study. In *Proceedings of the 2024 Joint International Conference on Computational Linguistics, Language Resources and Evaluation (LREC-COLING 2024)*, pages 8008–8015.

Ahmad Mustapha, Hadi Al-Khansa, Hadi Al-Mubasher, Aya Mourad, Ranam Hamoud, Hasan El-Husseini, Marwah Al-Sakkaf, and Mariette Awad. 2024. Arastem: A native arabic multiple choice question benchmark for evaluating llms knowledge in stem subjects. *arXiv preprint arXiv:2501.00559*.

Junho Myung, Nayeon Lee, Yi Zhou, Jiho Jin, Rifki Putri, Dimosthenis Antypas, Hsuvas Borkakoty, Eunsu Kim, Carla Perez-Almendros, Abinew Ali Ayele, et al. 2024. Blend: A benchmark for llms on everyday knowledge in diverse cultures and languages. *Advances in Neural Information Processing Systems*, 37:78104–78146.

Omer Nacar, Serry Taiseer Sibaee, Samar Ahmed, Safa Ben Atitallah, Adel Ammar, Yasser Alhabashi, Abdulrahman S Al-Batati, Arwa Alsehibani, Nour Qandos, Omar Elshehy, et al. 2025. Towards inclusive arabic llms: A culturally aligned benchmark in arabic large language model evaluation. In *Proceedings of the First Workshop on Language Models for Low-Resource Languages*, pages 387–401.

Abdulsalam obaid Alharbi, Abdullah Alsuhaibani, Abdulrahman Abdullah Alalawi, Usman Naseem, Shoaib Jameel, Salil Kanhere, and Imran Razzak. 2025. Evaluating large language models on health-related claims across arabic dialects. In *Proceedings of the 1st Workshop on NLP for Languages Using Arabic Script*, pages 95–103.

Zhaozhi Qian, Faroq Altam, Muhammad Alqurishi, and Riad Souissi. 2024. Cameleval: Advancing culturally aligned arabic language models and benchmarks. *arXiv preprint arXiv:2409.12623*.

Haneh Rhel and Dmitri Roussinov. 2025. Large language models and arabic content: a review. *arXiv preprint arXiv:2505.08004*.

Yuemei Xu, Ling Hu, Jiayi Zhao, Zihan Qiu, Kexin Xu, Yuqi Ye, and Hanwen Gu. 2025. A survey on multilingual large language models: Corpora, alignment, and bias. *Frontiers of Computer Science*, 19(11):1911362.

Sane Yagi, Ashraf Elnagar, and Shehdeh Fareh. 2023. A benchmark for evaluating arabic word embedding models. *Natural Language Engineering*, 29(4):978–1003.

Sane Yagi, Shehdeh Fareh, Ashraf Elnagar, Mariam Balajeed, Abdalla El-mneizel, and Mohammad Al-Badawi. 2024. Is arabic punctuation rule-governed? *Cogent Arts & Humanities*, 11(1):2303818.

Wenli Yang, Lilian Some, Michael Bain, and Byeong Kang. 2025. A comprehensive survey on integrating large language models with knowledge-based methods. *Knowledge-Based Systems*.

# A  Overview of Arabic LLM Benchmarks

Table 1: Overview of Arabic LLM Benchmarks

| Benchmark | Year | LLMs Evaluated | Task(s) | Dataset(s) Description |
|---|---|---|---|---|
| **AraT5** (Elmadany et al., 2022) | 2022 | AraT5 (Small, Base, Large, XL), mT5, mBART, AraGPT2, MAR-BERT | 8 tasks: Text Classification, Sentiment Analysis, Named Entity Recognition (NER), Extractive Question Answering (QA), Paraphrasing, Summarization, Headline Generation, Text Simplification | Data collected from 8 diverse Arabic sources including Arabic Wikipedia, OSCAR, OPUS, Tashkeela, SLSA, and others; resulting in a corpus of 200M sentences (50GB); preprocessed into a text-to-text format. |
| **Beyond English** (Kwon et al., 2023) | 2023 | GPT-4, ChatGPT-3.5 Turbo, LLaMA-7B, Vicuna-13B, Bactrian-Xbloom-7B, Bactrian-Xllama-7B | Grammatical Error Correction (GEC) | QALB-2014 (L1), QALB-2015 (L1 & L2) |
| **Dolphin** (Elmadany et al., 2023) | 2023 | Falcon-40B-Instruct, Falcon-180B-Chat, GPT-3.5-Turbo, GPT-4, ChatGPT | 10 NLG tasks: dialogue generation, question answering, data-to-text, storytelling, summarization, translation, paraphrasing, definition generation, classification, correction/refinement (includes code-switching and Arabizi) | 20K Arabic prompts with 200K completions, covering diverse topics and language forms, including Modern Standard Arabic, dialects, code-switched inputs, and Arabizi; prompts created by native speakers and aligned with high-quality completions |
| **Evaluating ChatGPT and Bard AI on Arabic Sentiment Analysis** (Al-Thubaity et al., 2023) | 2023 | GPT-3.5, GPT-4, Bard AI (PaLM 2) | Sentiment Analysis (Classification & Generation) | Saudi Dialect Twitter Corpus (SDTC): 2,690 used (558 positive, 1,632 negative, 500 neutral) |
| **Evaluation of Bard and ChatGPT on MT** (Kadaoui et al., 2023) | 2023 | ChatGPT (GPT-3.5-turbo), Bard | Machine Translation | 1,000 sentences from 10 Arabic dialects (100 per dialect) from the MADAR corpus, with corresponding MSA and English translations |
| **GPTAraEval** (Khondaker et al., 2023) | 2023 | ChatGPT-3.5, ChatGPT-4 | Text classification, natural language inference (NLI), question answering (QA), paraphrase identification, sentiment analysis, named entity recognition (NER), topic classification, hate speech detection, offensive language detection, dialect identification, translation, coreference resolution, headline generation, text summarization | 60 Arabic datasets covering Modern Standard Arabic and multiple Arabic dialects; varying in domain, size, and complexity; formatted for zero-shot prompt-based evaluation |
| **AlGhafa** (Almazrouei et al., 2023) | 2023 | AraT5, CAMeLBERT, mBERT, mGPT, GPT-3.5-turbo, AraGPT2-Mega, Noor-10B, Jais-13B, lGhafa-1B/3B/7B/14B | 45 tasks across 4 categories: knowledge, reasoning, reading comprehension, math & coding | 7,226 multiple-choice questions from diverse Arabic sources across linguistic and domain topics |
| **A Benchmark Evaluation of Multilingual LLMs for Arabic Cross-Lingual NER** (Al-Duwais et al., 2024) | 2024 | mBERT, XLM-R, BERTIN, ByT5, BLOOM, mT0 | NER | 7 Arabic NER datasets: ANERcorp, AQMAR, CAMeL, WikiFANE, Winerz, Arman, Arap-Tweet; domains: news, Wikipedia, social media |
| **ArabicMMLU** (Koto et al., 2024) | 2024 | GPT-3.5, GPT-4, BLOOMZ, mT0, LLaMA2, Falcon, XGLM, AraT5, AraGPT2, AceGPT, Jais (total 35 models) | Knowledge tasks | 14,575 Arabic multiple-choice questions from school exams in 8 Arabic-speaking countries; |
| **ArabLegalEval** (Hijazi et al., 2024) | 2024 | GPT-4, GPT-4o, GPT-3.5, Claude-3 Opus, Command R, Command R Plus, Llama3 (8B & 70B), Aya-101, Jais | MCQs, Open-ended QA, LegalBench QA (Consumer Contracts, Contracts, Privacy QA/Entailment) | 10,583 Arabic MCQs (from MoJ & BoE), 492 Najiz FAQs, 15,804 translated LegalBench samples, ArabicMMLU subset for legal reasoning benchmarking. |
| **ARADICE** (Mousi et al., 2025) | 2025 | Jais-13B, AceGPT-13B, Llama-3-8B, Mistral-7B, Fanar-8.7B, Qwen2.5-7B, Gemma2-9B, Aya-8B | Dialect Identification, Dialect Generation, Machine Translation, Commonsense Reasoning, World Knowledge, Reading Comprehension, Misinformation Detection, Cultural Understanding | 45K+ post-edited examples across QADI, ADI, ADD, MADAR, ArabicMMLU, PIQA, OBQA, Winogrande, BoolQ, Belebele, TruthfulQA, and AraDiCE-Culture |
| **AraSTEM** (Mustapha et al., 2024) | 2024 | AraT5, AraGPT2, MT0 (Small, Base, Large), XGLM (1.7B–7.5B), Bloomz (560M–7B1), AceGPT (7B, 13B), LLaMA (2 & 3.1), Falcon (7B, 40B), Jais (13B, 30B) | Zero-shot multiple-choice answering | 11,637 Arabic MCQs in STEM (math, biology, physics, IT, chemistry, pharmacy, medicine, dentistry); levels: primary, secondary, college; sourced via scraping, manual extraction, OCR from PDFs; annotated with source traceability |
| **AraTrust** (Alghamdi et al., 2025) | 2025 | GPT-3.5 Turbo, GPT-4, AceGPT 7B, AceGPT 13B, Jais 13B | Trustworthiness evaluation | 522 multiple-choice questions across 8 categories (Truthfulness, Ethics, Physical Health, Mental Health, Unfairness, Illegal Activities, Privacy, Offensive Language) and 34 subcategories, all human-written |
| **Halwasa** (Mubarak et al., 2024) | 2024 | GPT-3.5, GPT-4 | Factual sentence generation to evaluate models' hallucinations | 10K Arabic sentences (5K/model) generated using 1,000 random words from the SAMER corpus, annotated for factuality, correctness, linguistic errors, and references |
| **LARaBench** (Abdelali et al., 2023) | 2023 | GPT-3.5-turbo, GPT-4, BLOOMZ, Jais-13b-chat, Whisper, USM | 33 tasks across NLP and Speech | 61 publicly available datasets; 296K samples; 46h speech; 30 TTS sentences; covers MSA and dialects, across genres like news, tweets, telephony |

Table 1 – continued from previous page

| Benchmark | Year | LLMs Evaluated | Tasks | Dataset(s) Description |
|---|---|---|---|---|
| **Arabic Word/Location Sense Disambiguation** (Abdel-Salam, 2024) | 2024 | LLama3, LLama3-Instruct, WizardLM-2, AceGPT-7B, OpenChat | Word Sense Disambiguation (WSD), Location Mention Disambiguation (LMD) | WSD: SALMA corpus (100 train, 1,340 test); LMD: IDRISI-D (2,170 train, 333 val, 791 test) |
| **The Qiyas Benchmark** (Al-Khalifa and Al-Khalifa, 2024) | 2024 | ChatGPT-3.5-turbo, ChatGPT-4, Gemini-pro (partial) | Mathematical reasoning and Language understanding | 2,407 multiple-choice questions derived from Saudi Arabia's Qiyas GAT. Includes math, geometry, algebra, statistics, and five types of verbal tasks |
| **Jawaher** (Magdy et al., 2025) | 2025 | Llama-3.1-8B-Instruct, Llama-3.2-3B-Instruct, Gemma-2-9B-it, GPT-4o, Gemini 1.5 Pro, Claude 3.5 Sonnet, Cohere Command R+ | Translation, Explanation | 10,037 Arabic proverbs from 20 dialects with idiomatic/literal English translations, Arabic and English explanations. |
| **ArabicSense** (Lamsiyah et al., 2025) | 2025 | Gemma, LLaMA-3, Mistral-7b | Commonsense Validation, Multiple-Choice Explanation, Generative Explanation | 3954 train, 848 val, 848 test samples per task from Arabic Wikipedia, generated using GPT-4 |
| **Arabic Fact-Checking** (Gupta et al., 2025) | 2025 | Llama 3 8B, Llama 3 70B, GPT-3.5-turbo, Gemini 1.0 Pro | Arabic fact-checking (binary classification: true/false) | 771 Arabic claims from X-Fact dataset (filtered for 'true' or 'false' only; 730 false, 41 true) |
| **Health-Related Claims Across Arabic Dialects** (obaid Alharbi et al., 2025) | 2025 | GPT-4 | Health claim verification across dialects | 329 claims (191 from AraFacts + 138 from ArCOV19-Rumors), categorized as true, false, mixed |
| **Evaluation of LLMs on Arabic Punctuation Prediction** (Al Wazrah et al., 2025) | 2025 | GPT4-o, Gemini 1.5, JAIS-13B, AceGPT-13B, SILMA-9B, ALLaM-1, CommandR+, AraBERT | Arabic punctuation prediction | 10,046 annotated Arabic paragraphs from 25 books, manually cleaned and tokenized, covering six punctuation marks; split into training, validation, and test sets |
| **HalluVerse25** (Abdaljalil et al., 2025) | 2025 | GPT-4o, GPT-4o-mini, phi-4, PaLM 2, Mistral-7b, Qwen-2.5 (7b, 72b), LLaMA-3.3, Gemini, Gemma | Hallucination Detection | 3116 factual + hallucinated pairs (biography-based) in English, Arabic, Turkish |
| **Palm** (Alwajih et al., 2025) | 2025 | GPT-4o, Claude-3.5-Sonnet, Command R+ (104B), Qwen2.5-72B, Qwen2.5-7B, Qwen2.5-3B, Qwen2.5-1.5B, JAIS-13B, AceGPT-v2-32B, AceGPT-v2-8B, LLaMA-3.1-70B, LLaMA-3.1-8B, LLaMA-3.2-3B, LLaMA-3.2-1B, Gemma-2-27B, Gemma-2-9B, Gemma-2-2B, Phi-3.5-mini (18 models) | To benchmark LLMs' capabilities in culturally-aware and dialect-specific instruction following across the Arab world. It evaluates LLMs' ability to understand and generate culturally relevant, linguistically appropriate responses in Arabic dialects and MSA. | 17,411 human-authored Arabic instruction–response pairs (MSA and 10 dialects) across 22 Arab countries and 20 cultural domains; includes train, public test, and private test splits |
| **SaudiCulture** (Ayash et al., 2025) | 2025 | GPT-4, Llama 3.3, FANAR, Jais, AceGPT | Cultural understanding, QA (open-ended, single-answer, and multi-answer formats) | SaudiCulture: 441 questions across 5 Saudi regions + general, covering 8 cultural domains (food, clothing, celebrations, etc.) in open-ended, single-answer, and multi-answer formats |
| **Towards Inclusive Arabic LLMs** (Nacar et al., 2025) | 2025 | Qwen2.5-72B-Instruct, CohereForAI/aya-expanse-32b, Qwen2.5-32B-Instruct, Google/Gemma-2, SILMA-9B, FreedomIntelligence/AceGPT, JAIS-family, LLaMA models | Multitask Language Understanding | Refined Arabic MMLU benchmark with over 14,000 questions, including 2,466 culturally sensitive questions and 766 culturally enriched additions (e.g., Islamic Religion, Islamic Ethics, Old Arab History). |

# Beyond the Boundaries of Research Fields - Mapping Educational Science in the Broader Academic Discourse about Artificial Intelligence

**Martin Rhem, Zhiru Sun, Maxime Holmberg Sainte-Marie**
Department of Design, Media & Educational Science, University of Southern Denmark
{mareh,zhiru, mhsm}@sdu.dk

## Abstract

Over the past three decades, research in Artificial Intelligence (AI) has steadily grown and advanced. The recent surge in generative AI tools, such as ChatGPT, has reignited interest also among research fields that otherwise seemingly turned their focus elsewhere. More specifically, the field of Education has seen a remarkable surge in research on AI (AIEd), examining its impact on various types of learning. While these investigations offer valuable exploratory insights, many lack systematic analysis. Moreover, they overlook the potential overlap and commonalities with other research domains. This study addresses this gap by utilizing the OpenAlex database to identify and analyze not only AI-related publications since 2000. Through part-of-speech tagging and semantic network analysis, and focusing on the AIEd as our anchor, we map content relationships across the literature to reveal thematic structures and potential synergies across research fields. The findings offer an exploratory overview of the evolving research landscape, the position of AIEd therein, and suggest directions for future inquiry.

## 1 Introduction

There has been a growing interest in AIEd research over the past three decades, with a focus on three key areas: predictions, personalization, and assessment. In the area of prediction, AI tools play a pivotal role in learning analytics by analyzing learner profiles to predict academic achievement, dropout risks, retention rates, and even admission decisions. Numerous studies leverage these insights to enable institutions to intervene proactively, support at-risk students, and make data-driven decisions to enhance educational outcomes (Batool et al., 2023; Fahd et al., 2022). Tools such as Course Signals, Coursera, and learning management systems exemplify the use of predictive analytics, applying AI to monitor students' learning process and enhance their academic success. For personalization,

AI analyzes student data to create tailored learning pathways that address individual needs and optimize outcomes. Intelligent tutoring systems, such as Carnegie Learning's Cognitive Tutor, offer personalized guidance and feedback, while adaptive learning platforms like Khan Academy and Duolingo adjust content dynamically based on student progress. These tools not only foster deeper engagement but also empower students to achieve their academic goals at their own pace (Lin et al., 2023). In assessment, AI tools have been used to automate grading and assessment tasks, allowing educators to focus on providing meaningful, targeted feedback. Tools like Turnitin's automated essay scoring system efficiently evaluate written work, saving instructor's time for personalized feedback to students (Kostka and Toncelli, 2023). Together, the AI tools applied in these three key areas primarily fall under reactive and limited memory AI. However, the rapid advancement of generative AI (GAI) tools has opened up new possibilities. Starting in 2022, there has been a surge in AIEd research, with a focus on GAI tools like ChatGPT and their transformative impact on teaching and learning.

## 2 Rationale of the Study

Researchers are exploring how GAI tools can be applied in various research fields, including computer science (Kar et al., 2023; Parker, 2025), health sciences (Moulaei et al., 2024; Sai et al., 2024a), engineering (Sai et al., 2024b; Vu et al., 2024), as well as business and economics (Orchard and Tasiemski, 2023; Yu and Qi, 2024). More specifically, and from the perspective of AIEd, research is largely dealing with highlighting its potential to reshape traditional educational practices.

This study builds on the considerations laid out in the previous section and investigates the position of AIEd in the broader academic discourse about AI. More specifically, we employ part-of-

speech tagging and semantic network analyses to analyze publications on AI identified in the OpenAlex database. This type of approach has been suggested to map research fields in a wider context (McAllister et al., 2022) and highlight the potential for synergies and spillovers (Hou et al., 2022). The goal is to unveil underlying structures that can be used as a point of departure for further investigations of common terminology, content topics, and interrelations between research fields. Consequently, in the context of this exploratory study, our research questions are:

RQ1) *What does the general academic landscape look like on the overarching research topic of AI?*

RQ2) *To what extent can we identify content overlaps between research fields?*

RQ3) *Where does AIEd interconnect with other research fields?*

Next, we provide an overview of our data collection procedures and two types of analysis to guide the reader.

## 3 Methods

### 3.1 Data Collection

Article metadata was extracted from OpenAlex, the first open-source, large-scale, and multilingual bibliometric database (Priem et al., 2022). Based mainly on data extracted from the discontinued Microsoft Academic Graph as well as from other bibliometric and bibliographic databases and repositories, OpenAlex has been shown to offer superior document, journal, and language coverage than existing proprietary databases (Alperin et al., 2024; Culbert et al., 2025; Jiao et al., 2023; Thelwall and Jiang, 2025). As a result, a growing number of research projects are using this database to conduct various types of bibliometric-based research, focusing on topics such as article retractions (Hauschke and Nazarovets, 2025; Ortega and Delgado-Quirós, 2024; Yiru et al., 2025), science mapping (Haunschild and Bornmann, 2024), open access publishing (Simard et al., 2025), data reuse (Krause and Mongeon, 2023), as well as quality of geographic, disciplinary, or linguistic coverage (Maddi et al., 2025; Céspedes et al., 2025).

Of particular interest and relevance here is the recently implemented article-based topic classification of articles indexed in the database (Barrett, 2024). Based on proven multilevel coarsening and refinement article classification procedure developed and used at CWTS (Eck and Waltman, 2024; Waltman and Van Eck, 2012), these topic categories are increasingly used by geographically and disciplinary diverse scientific communities for various research purposes (Arroyo-Machado and Costas, 2023; Cebrián et al., 2025; Couto and Baltazar, 2025; de Carvalho Segundo et al., 2024)

Our search terms were deliberately broad to cast a wide net in the search. Search terms included a mix of "artificial intelligence" and specific terms that are generally mentioned in literature reviews on the topic across different disciplines, including "risk", "challenges", "opportunities", "education", and "impact". While a "casting the net wide" approach also has potential drawbacks, we believe that this technique provides a valuable approach to possibly discover commonalities between disciplines (Authors, 2019). The search was conducted on the 7th of May, 2025, and was limited to publications ranging from the 1st of January 2000 until the 31st of December 2024. The 43,598 results were downloaded into a dataframe in the statistical software package R. We then filtered for publications that also included an abstract, resulting in a consolidated dataset of 27,202, which were then analyzed using the R libraries igraph, quanteda, tm, and udpipe.

### 3.2 Part-of-Speech Tagging

The main idea of part-of-speech tagging is to assign each word of a text to its proper syntactic tag in the context of its appearance (Chiche and Yitagesu, 2022). This is also referred to as grammatical tagging (Khan et al., 2019) and ensures grammatical relevance. Here, we used POS on the titles and abstracts of the collected publications. More specifically, we focused on nouns, verbs and adjectives. Furthermore, based on the POS results, we then determined n-grams to capture meaningful co-occurrence patterns (Bai et al., 2021; Ojo et al., 2021). This complementary approach allowed us to identify key concepts and their contextual relationships from both titles and abstracts.

### 3.3 Semantic Network Analysis

Semantic network analysis (SemNA)is a technique used to identify and visualize relationships between key concepts based on their co-occurrence within textual data (Castelblanco et al., 2021; Segev, 2022). In the context of this study, key concepts were determined by the classification procedure

employed by OpenAlex. We specifically focused on the category of "subfield", as it allowed us to take a more granular approach to looking at the data. The textual data was taken from the preceding POS analyses and determined n-grams. Here, we applied SemNA to both titles and abstracts of the collected publications. Titles offered a high-level overview of thematic structures, while abstracts allowed for a deeper exploration of concept relationships. The resulting networks then provide a more nuanced understanding of the topical landscape and its underlying structures.

## 4 Results

### 4.1 Descriptive Analysis

Tables 1 and 2 provide a glimpse at the underlying data for the categories fields and subfields.

Table 1: Top Research Fields (based on N publications in our data set)

| Field | N |
|---|---|
| Computer Science | 16687 |
| Medicine | 7819 |
| Social Sciences | 5753 |
| Engineering | 3226 |
| Business, Management & Accounting | 2525 |

Table 2: Top Research Subfields (based on N publications in our data set)

| Subfield | N |
|---|---|
| Artificial Intelligence | 8463 |
| Health Informatics | 4293 |
| Computer Science Applications | 3449 |
| Safety Research | 3011 |
| Information Systems | 2422 |
| Management Information Systems | 1351 |
| Radiology, Nuclear Medicine & Imaging | 1250 |
| Management Science & Operations Research | 948 |
| Education | 823 |

### 4.2 POS, n-grams & SemNA

Starting with the investigation of titles, and focusing on the subfields indicated in Table 2, we determined the semantic network visualized in Figure 1.

Here, we see that there is considerable content overlap between the different research sub-



Figure 1: Semantic Network for Subfields (Titles)



Figure 2: Semantic Network for AIEd (Titles)

fields. More specifically, irrespective of the underlying discipline, scholars were largely concerned with "explainable AI", "risk management", "challenges", and "case studies". Focusing on AIEd, Figure 2 underlines these findings, while clearly showing that the perspective on these topics has been from an educational perspective, e.g. concentrating on the context of higher education and generative tools. The later suggests an inerest and discourse about how

We then turned to the abstracts, in order to engage into a deeper exploration of concept relationship. Figure 3 provides the overall view, while Figure 4 highlights AIEd again. Here, we found that the core of the semantic network is largely driven by the subfields of "artificial intelligence", "information systems", "management science", and "industrial and manufacturing engineering". Another closer look at AIEd (Figure 4) then revealed that scholars seem particularly interested in AI for "language education", "student engagement", and "prompt engineering". Moreover, a sizeable amount of publications were also concerned with "compliance" and "risk management".

Figure 3: Semantic Network for Subfields (Abstracts)



Figure 4: Semantic Network for AIEd (Abstracts)

## 5 Discussion

This study set out to map AIEd in the broader academic discourse about AI and identify possible synergies and spillovers across research fields. Our results indicate that while AIEd has certainly experienced a surge in publications, other traditionally more technical and computational domains remain at the core of the research field at this point in time (RQ1). However, using POS, n-grams and semantic networks, we have also been able to show considerable interconnections between the different fields (RQ2). This supports the general trend of more interdisciplinary research and the specific request to address the topic of AI from different, interrelated perspectives (Følstad et al., 2021; Newman, 2024). Finally, concentrating on AIEd, we have been able to identify promising "low-hanging fruits" of content overlaps with different research fields that can provide mutual benefit for AIEd and related fields (RQ3). More specifically, while AIEd can benefit from groundwork and various applications on the topic of e.g. "prompt engineering", other fields could use AIEd as a type of proving ground to test and validate the findings from their own fields. While this study provides valuable insights into where AIEd is located in the larger context of research on AI, it is also subject to some limitations

that should be considered when interpreting the results and designing future studies on the topic. First, we built our literature search on topics and concepts informed by discourses within the research field of education. Future research, also based on or findings should consider casting the web even wider, in order to potentially collect a more holistic sample of the underlying academic discourse. Second, the focus of our work has been on titles and abstracts, rather than full papers. While this is a good point of departure, it bases the analyses on limited amount of text. Next possible steps could include the analyses of entire publications, to provide an even more nuanced view, providing even better chances to identify possible overlaps and synergies.

## References

Juan Pablo Alperin, Jason Portenoy, Kyle Demes, Vincent Larivière, and Stefanie Haustein. 2024. An analysis of the suitability of openalex for bibliometric analyses. *arXiv preprint arXiv:2404.17663*.

Wenceslao Arroyo-Machado and Rodrigo Costas. 2023. Do popular research topics attract the most social attention? a first proposal based on openalex and wikipedia. In *27th International Conference on Science, Technology and Innovation Indicators (STI 2023)*. International Conference on Science, Technology and Innovation Indicators.

Huiwen Bai, Guangjie Liu, Weiwei Liu, Yingxue Quan, and Shuhua Huang. 2021. N-Gram, Semantic-Based Neural Network for Mobile Malware Network Traffic Detection. *Security and Communication Networks*, 2021:1–17.

J.P. Barrett. 2024. Open alex topic classification. *GitHub repository*, https://github.com/ourresearch/openalex-topic-classification.

Saba Batool, Junaid Rashid, Muhammad Wasif Nisar, Jungeun Kim, Hyuk-Yoon Kwon, and Amir Hussain. 2023. Educational data mining to predict students' academic performance: A survey study. *Education and Information Technologies*, 28(1):905–971.

Gabriel Castelblanco, Jose Guevara, Harrison Mesa, and Andres Sanchez. 2021. Semantic Network Analysis of Literature on Public-Private Partnerships. *Journal of Construction Engineering and Management*, 147(5):04021033.

Guillem Cebrián, Ángel Borrego, and Ernest Abadal. 2025. Openalex y crossref como fuentes de datos bibliográficas alternativas a web of science y scopus en ciencias de la salud. *Revista Española de Documentación Científica*, 48(1):1649–1649.

Lucía Céspedes, Diego Kozlowski, Carolina Pradier, Maxime Holmberg Sainte-Marie, Natsumi Solange Shokida, Pierre Benz, Constance Poitras, Anton Boudreau Ninkov, Saeideh Ebrahimy, Philips Ayeni, et al. 2025. Evaluating the linguistic coverage of openalex: An assessment of metadata accuracy and completeness. *Journal of the Association for Information Science and Technology*.

Alebachew Chiche and Betselot Yitagesu. 2022. Part of speech tagging: a systematic review of deep learning and machine learning approaches. *Journal of Big Data*, 9(1):1–25. Publisher: SpringerOpen.

João Couto and Maria Emilia Baltazar. 2025. Sustainable airport development: A literature review based on preferred reporting items for systematic reviews and meta-analyses methodology, using openalex database. *Sustainability*, 17(9):4184.

Jack H Culbert, Anne Hobert, Najko Jahn, Nick Haupka, Marion Schmidt, Paul Donner, and Philipp Mayr. 2025. Reference coverage analysis of openalex compared to web of science and scopus. *Scientometrics*, 130(4):2475–2492.

Washington Luís Ribeiro de Carvalho Segundo, Fabio Lorensi do Canto, Adilson Luiz Pinto, and Daniel Lima Sundfeld. 2024. Atração entre periódicos brasileiros de medicina: análise a partir de dados de citação do openalex. *Encontro Brasileiro de Bibliometria e Cientometria*, 9:1–7.

Nees Jan van Eck and Ludo Waltman. 2024. An open approach for classifying research publications.

Kiran Fahd, Sitalakshmi Venkatraman, Shah J Miah, and Khandakar Ahmed. 2022. Application of machine learning in higher education to assess student academic performance, at-risk, and attrition: A meta-analysis of literature. *Education and Information Technologies*, pages 1–33.

Asbjørn Følstad, Theo Araujo, Effie Lai-Chong Law, Petter Bae Brandtzaeg, Symeon Papadopoulos, Lea Reis, Marcos Baez, Guy Laban, Patrick McAllister, Carolin Ischen, Rebecca Wald, Fabio Catania, Raphael Meyer Von Wolff, Sebastian Hobert, and Ewa Luger. 2021. Future directions for chatbot research: an interdisciplinary research agenda. *Computing*, 103(12):2915–2942.

Robin Haunschild and Lutz Bornmann. 2024. The use of openalex to produce meaningful bibliometric global overlay maps of science on the individual, institutional, and national levels. *PloS one*, 19(12):e0308041.

Christian Hauschke and Serhii Nazarovets. 2025. (non-) retracted academic papers in openalex. *Journal of Information Science*, page 01655515251322478.

Xiaojing Hou, Ruichang Li, and Zhiping Song. 2022. A Bibliometric Analysis of Wicked Problems: From Single Discipline to Transdisciplinarity. *Fudan Journal of the Humanities and Social Sciences*, 15(3):299–329.

Chenyue Jiao, Kai Li, and Zhichao Fang. 2023. How are exclusively data journals indexed in major scholarly databases? an examination of the web of science, scopus, dimensions, and openalex. *arXiv preprint arXiv:2307.09704*.

Arpan Kumar Kar, P. S. Varsha, and Shivakami Rajan. 2023. Unravelling the Impact of Generative Artificial Intelligence (GAI) in Industrial Applications: A Review of Scientific and Grey Literature. *Global Journal of Flexible Systems Management*, 24(4):659–689.

Wahab Khan, Ali Daud, Khairullah Khan, Jamal Abdul Nasir, Mohammed Basheri, Naif Aljohani, and Fahd Saleh Alotaibi. 2019. Part of speech tagging in urdu: Comparison of machine and deep learning approaches. *IEEE Access*, 7:38918–38936. Publisher: IEEE.

Ilka Kostka and Rachel Toncelli. 2023. Exploring applications of chatgpt to english language teaching: Opportunities, challenges, and recommendations. *Tesl-Ej*, 27(3):n3.

Geoff Krause and Philippe Mongeon. 2023. Measuring data re-use through dataset citations in openalex. In *27th International Conference on Science, Technology and Innovation Indicators (STI 2023)*. International Conference on Science, Technology and Innovation Indicators.

Chien-Chang Lin, Anna YQ Huang, and Owen HT Lu. 2023. Artificial intelligence in intelligent tutoring systems toward sustainable education: a systematic review. *Smart Learning Environments*, 10(1):41.

Abdelghani Maddi, Marion Maisonobe, and Chérifa Boukacem-Zeghmouri. 2025. Geographical and disciplinary coverage of open access journals: Openalex, scopus, and wos. *PloS one*, 20(4):e0320347.

James T. McAllister, Lora Lennertz, and Zayuris Atencio Mojica. 2022. Mapping A Discipline: A Guide to Using VOSviewer for Bibliometric and Visual Analysis. *Science & Technology Libraries*, 41(3):319–348.

Khadijeh Moulaei, Atiye Yadegari, Mahdi Baharestani, Shayan Farzanbakhsh, Babak Sabet, and Mohammad Reza Afrash. 2024. Generative artificial intelligence in healthcare: A scoping review on benefits, challenges and applications. *International Journal of Medical Informatics*, page 105474. Publisher: Elsevier.

Joshua Newman. 2024. Promoting Interdisciplinary Research Collaboration: A Systematic Review, a Critical Literature Review, and a Pathway Forward. *Social Epistemology*, 38(2):135–151.

O. E. Ojo, A. Gelbukh, H. Calvo, and O. O. Adebanji. 2021. Performance study of n-grams in the analysis of sentiments. *Journal of the Nigerian Society of Physical Sciences*, pages 477–483.

Tim Orchard and Leszek Tasiemski. 2023. The rise of generative AI and possible effects on the economy. *Economics and business review*, 9(2):9–26. Publisher: Poznan University of Economics.

José Luis Ortega and Lorena Delgado-Quirós. 2024. The indexation of retracted literature in seven principal scholarly databases: a coverage comparison of dimensions, openalex, pubmed, scilit, scopus, the lens and web of science. *Scientometrics*, 129(7):3769–3785.

Jessica Parker. 2025. Generative AI (GAI) Use for Cybersecurity Resilience: A Scoping Literature Review. *International Journal of Applied Science*, 8(2):p1–p1.

Jason Priem, Heather Piwowar, and Richard Orr. 2022. Openalex: A fully-open index of scholarly works, authors, venues, institutions, and concepts. *arXiv preprint arXiv:2205.01833*.

Siva Sai, Aanchal Gaur, Revant Sai, Vinay Chamola, Mohsen Guizani, and Joel JPC Rodrigues. 2024a. Generative ai for transformative healthcare: A comprehensive study of emerging models, applications, case studies and limitations. *IEEE Access*. Publisher: IEEE.

Siva Sai, Revant Sai, and Vinay Chamola. 2024b. Generative AI for Industry 5.0: Analyzing the impact of ChatGPT, DALLE, and other models. *IEEE Open Journal of the Communications Society*. Publisher: IEEE.

Elad Segev. 2022. *Semantic network analysis in social sciences*. Routledge London.

Marc-André Simard, Isabel Basson, Madelaine Hare, Vincent Larivière, and Philippe Mongeon. 2025. Examining the geographic and linguistic coverage of gold and diamond open access journals in openalex, scopus and web of science. *Quantitative Science Studies*, pages 1–29.

Mike Thelwall and Xiaorui Jiang. 2025. Is openalex suitable for research quality evaluation and which citation indicator is best? *arXiv preprint arXiv:2502.18427*.

Thai-Hoc Vu, Senthil Kumar Jagatheesaperumal, Minh-Duong Nguyen, Nguyen Van Huynh, Sunghwan Kim, and Quoc-Viet Pham. 2024. Applications of generative AI (GAI) for mobile and wireless networking: A survey. *IEEE Internet of Things Journal*. Publisher: IEEE.

Ludo Waltman and Nees Jan Van Eck. 2012. A new methodology for constructing a publication-level classification system of science. *Journal of the American Society for Information Science and Technology*, 63(12):2378–2392.

Liu Yiru, Liu Yi, and Yuan Zihan. 2025. A comprehensive bibliometric analysis of retracted chapters based on openalex database. *Scientometrics*, 130(4):2425–2444.

Jason Yu and Cheryl Qi. 2024. The impact of generative AI on employment and labor productivity. *Review of Business*, 44(1):53–67. Publisher: St. John's University.

# From Performance to Process: Temporal Information Dynamics in Language Model Fine-tuning

**Frida Hæstrup**[1,2] and **Ross Deans Kristensen-McLachlan**[2]

[1]Dept. of Affective Disorders, Aarhus University Hospital – Psychiatry, Aarhus, Denmark
[2]Center for Humanities Computing, Aarhus University, Aarhus, Denmark
frihae@clin.au.dk, rdkm@cc.au.dk

## Abstract

Large language model performance has advanced rapidly in recent years, driven by technical improvements in areas like model architecture, scaling, and reinforcement learning. However, much of our understanding of these models remains rooted in static evaluations calculated post-training. While informative, these snapshots offer limited insight into how models learn, adapt, and transform internally during training, overlooking dynamic processes and representational shifts that occur throughout fine-tuning, potentially concealing important aspects of model behavior. We aim to contribute to ongoing efforts to open the 'black box' of language models by analyzing temporal information dynamics during fine-tuning. Our findings suggest that tracking these internal dynamics demonstrates both training-regime-specific and task-specific differences in learning and may eventually contribute to applications such as change point detection or adaptive training strategies. Ultimately, this work moves toward a more nuanced, mathematical formulation of what learning does to a model, highlighting the constant flux of representational change that underlies seemingly stable performance improvements.

## 1 Introduction

The rapid development and widespread deployment of large language models (LLMs) have amplified interest in understanding how these models function internally. In pursuit of improved model performance and generalization, the development of pre-trained LLMs has led to models that are increasingly becoming larger and more complex (Simon, 2021; Brown et al., 2020). Such complexity, often driven by millions or even billions of parameters, enables these models to capture and learn intricate patterns within the training data, allowing them to achieve state-of-the-art results across a wide array of tasks (Devlin et al., 2019; Wang et al., 2018; Rozière et al., 2024; Wang et al., 2020).

However, this power comes at a significant cost: it obscures the internal mechanisms by which models arrive at their predictions, rendering the path from input to output difficult to interpret and explain. As these models are increasingly adopted in sensitive and high-stakes domains, the need for transparency into their internal processes becomes not just desirable, but essential (Hassija et al., 2024; Embarak, 2023; Chen et al., 2025). To better understand what these models are actually learning — and how their internal states evolve during training— we must look beyond static evaluations and examine the learning process itself. Standard evaluation metrics such as accuracy, perplexity, or F1-score provide only static snapshots of model behavior. These metrics reflect *what* a model achieves but offer little insight into *how* it learns.

In this paper, we propose analyzing the temporal dynamics of learning in language models within an information-theoretic framework (MacKay, 2002), conceptualizing a model's internal state as a dynamic information system (**?**). Rather than focusing solely on final performance, we track how internal representations evolve during fine-tuning. This allows us to characterize learning as a continuous sequence of representational shifts, offering a more granular and process-oriented perspective on model behavior.

Our proposed framework builds on a growing body of research that has used information theory to study the evolution of complex, dynamic systems. In particular, several studies have modeled cultural and linguistic phenomena by analyzing the balance between how much new information is being introduced and that information's longevity within the system (Barron et al., 2018; Nielbo et al., 2021a,b; Vrangbæk and Nielbo, 2021; Wevers et al., 2021; Krisensen-McLachlan et al., 2024). These studies used windowed relative entropy to quantify the **novelty** of a system - the extent to which a given time period diverges from preced-

ing time periods - and the **resonance** of a system, which captures how information persists over time.

We extend this framework to the context of deep learning by treating the internal states of a language model as a dynamic information system. We investigate the evolution of internal information structure in models from the English BERT family (Devlin et al., 2019) as they are fine-tuned across various classification tasks. Through a series of controlled experiments, we continuously extract internal representations from different BERT models throughout the fine-tuning process. We adopt an exploratory approach, examining whether tracking the dynamics of internal representations over time can reveal novel insights into the mechanisms of learning within these models.

We argue that this approach offers a rich perspective on what it means for a model to learn and opens the door to future applications, such as tracking learning trajectories, identifying shifts in representational focus, or detecting meaningful change points during training. Ultimately, we aim to bridge the gap between surface-level performance and deeper representational change, providing insight into the temporal structure of learning itself.

## 1.1 Related Work

Prior research has explored how fine-tuning affects the internal structure of transformer-based models such as BERT. A common approach involves probing internal layers to identify which aspects of the model change during adaptation to downstream tasks (Phang et al., 2021; Hao et al., 2020; Merchant et al., 2020; Zhou and Srikumar, 2022; Voita and Titov, 2020; Liu et al., 2019; Tenney et al., 2018; Voita and Titov, 2020). Hao et al. (2020) employ divergence-based measures to track shifts in attention patterns and find that fine-tuning primarily alters the attention modes of higher layers. This is consistent with observations from Merchant et al. (2020) who use probing classifiers and ablation experiments to show that representational change during fine-tuning is concentrated in upper layers. Furthermore, they find variations in this effect across fine-tuning tasks. For example, tasks such as dependency parsing produce deeper representational shifts than tasks like natural language inference or reading comprehension.

Further analyses have investigated the spatial structure of learned representations (Coenen et al.

(2019); Hernandez and Andreas (2021). Comparing the spatial structure of class-level embeddings before and after fine-tuning, Zhou and Srikumar (2022) observe that class representations are pushed further apart in the embedding space after fine-tuning, even in cases where the classes were already linearly separable. Extending the findings of Merchant et al. (2020), they also report that while higher layers change more than lower ones, these changes preserve structural similarity with the pre-trained model, suggesting that fine-tuning reshapes but does not fully overwrite earlier representations.

While these studies offer valuable insight into how models change across fine-tuning, they are typically limited to static comparisons between pre-trained and post-trained states. In contrast, our work adopts a dynamic perspective, examining internal representations at every step *during* the fine-tuning process. Moreover, rather than analyzing intermediate encoder layers, we focus on prediction-layer outputs, treating class-level output vectors as a dynamic system whose evolution reflects learning in real time. This allows us to capture transient changes and transitions that static snapshots may miss, offering a more granular view of representational dynamics during training.

## 2 Methods

We base our analysis on information signals extracted from 24 experiments: four pre-trained large language models fine-tuned on three classification tasks under two conditions. Details of this process are laid out in the following sections.[1]

### 2.1 Model architectures

We fine-tune four different pre-trained BERT-style models, namely **BERT** (Devlin et al., 2019), **distilBERT** (Sanh et al., 2020), **roBERTa** (Zhuang et al., 2021), and multilingual BERT (**mBERT**) (Devlin et al., 2019). The models are all core models that have been trained across many language-understanding tasks. Each model is based primarily on the BERT architecture, although they each display variations across different parameters such as size or training regime, allowing for a range of possible comparisons across models. An overview of the key differences across model types can be found in Appendix A.2. The pre-trained model weights of

---

[1]The code-base for the project can be found at https://github.com/frillecode/BERT-infodynamics

all four models were retrieved from HuggingFace.[2]

## 2.2 Classification tasks

We fine-tune the above-mentioned pre-trained models across three different language classification tasks from the General Language Understanding Evaluation (GLUE) benchmark (Wang et al., 2018). GLUE comprises a collection of resources for evaluating the performance of natural language understanding systems across a wide range of linguistic tasks. GLUE consists of nine different language understanding tasks, each built on established English-language text datasets, that are widely accepted as standard benchmarks for assessing how well models can understand and process natural language (Devlin et al., 2019; Radford et al., 2019). In the present study, a subset of three tasks from the GLUE benchmark is used, namely:

- **MNLI**: The Multi-Genre Natural Language Inference Corpus (Williams et al., 2018).

- **MRPC**: The Microsoft Research Paraphrase Corpus (Dolan and Brockett, 2005).

- **SST-2**: Stanford Sentiment Treebank (Socher et al., 2013).

The choice of using a subset of tasks is motivated by the following reasons. Firstly, the GLUE benchmark is typically used to assess how well models generalize across tasks and text genres, often with the ultimate goal of driving the development of robust natural language understanding systems (Wang et al., 2018). In contrast, this study seeks to explore the underlying processes of the models as they learn rather than assessing their final performance. Secondly, the experiments in this study make for 24 different fine-tuning processes and subsequent analyses, with the windowed relative entropy calculation adding substantial computational load. Thirdly, the choice of these tasks ensures that the study encompasses both binary and multi-class classification problems, as well as different dataset sizes. Furthermore, the tasks cover a wide range of linguistic phenomena as they represent each of the three general categories of the benchmark (Wang et al., 2018). As such, the tasks provide a sufficient variety of linguistic challenges to, within the scope

---

| Hyperparameter | Values |
|---|---|
| Batch size | $16, 32$ |
| Learning rate | $2e^{-5}, 3e^{-5}, 5e^{-5}$ |
| $N$ epochs | $2, 3, 4$ |

Table 1: Search space for hyperparameter optimization.

of this study, explore how models process and learn throughout different natural language understanding tasks.

## 2.3 Training procedures

The fine-tuning process of each model on each task is carried out under two conditions differing in parts of the training setup.

In the **fixed condition**, the hyperparameters of the training process are kept fixed across all experiments to allow for a more direct comparison. The models are trained for 5000 steps using a batch size of 64. All other hyperparameters are kept at default values.

In the **optimized condition**, hyperparameter optimization is incorporated in the training process to explore the effects of optimizing the models' learning process to the task. We perform a simple grid search over pre-defined values for batch size, learning rate, and number of epochs. We use the search space recommended in the original BERT paper (Devlin et al., 2019), as seen in Table 1. For each experiment, we run a total of 10 trials. Based on this, we define the best configuration for each experiment. These can be found in the Appendix A.1. All other hyperparameters are kept at default values.

In both conditions, a standard pipeline for fine-tuning machine-learning models was employed using the HuggingFace Transformers library (v. 4.42.4) (Wolf et al., 2020), and the datasets for the different GLUE tasks were retrieved using the Datasets class. All models are fine-tuned using a Cross-Entropy loss function, and the standard training and validation splits are retrieved automatically upon accessing the datasets from GLUE. All analysis is performed using Python (v. 3.12.3).

## 2.4 Feature extraction

During the fine-tuning process for each of the experiments, we save the logits at every training step by extracting the output of the last layer of the neural network. By passing the logits through the softmax function, they are converted to vectors rep-

Figure 1: Overview of pipeline for extracting information signals from logits. Fictional example for a classification task with 3 classes and a window size of $w = 2$. **(a)** illustrates the matrix with probability scores based on which the information signals will be extracted. **(b)** illustrates a novelty signal, with the blue area representing the window size within which it is calculated and the light red area representing the documents that are removed. **(c)** illustrates a transience signal, with the green area representing the window size within which it is calculated and the red area representing the documents that are removed.

resenting a probability distribution across labels. As the fine-tuning process continues over training steps, the resulting matrix becomes a temporally sorted series of probability distributions representing the model's predictions. Since these matrices (one for each experiment) capture how the models' predictions evolve over time, this can be used as a proxy reflecting the learning process as the models update their internal representations in response to the data. These probability distribution matrices hence serve as the input from which to extract information signals, as described in the following section. A visual representation of the process of extracting the information signals from the logits can be seen in Figure 1.

## 2.5 Information dynamics

Based on the temporally sorted probability scores for each of the experiments, we employ methods from information theory to extract information sig-

nals (novelty, resonance, transience). Using windowed relative entropy, we can measure the similarity (or 'surprise') between the information patterns in a series of probability distributions (Cover and Thomas, 2006). **Novelty** serves as a measure of how surprising the probability distribution patterns in a document are given past documents, **transience** measures the extent to which those patterns persist in future documents, and **resonance** measures the degree to which patterns in future documents conform to the novelty.

Information signals are extracted for each document using a window size of 160 ($w$=160). A document in this context refers to a document from the training data (i.e. an input sentence) of the given GLUE task that the model sees during fine-tuning. As such, a window size of 160 means that the information signals are extracted by comparing the model's representation of the current input sentence to the previous 160 input sentences and the

following 160 input sentences.

For the implementation of relative entropy, Jensen-Shannon divergence was used:

$$JSD(s^{(j)} \mid s^{(k)}) = \frac{1}{2}D(s^{(j)} \mid M) + \frac{1}{2}D(s^{(k)} \mid M)$$

(1)

where $M = \frac{1}{2}(s^{(j)} + s^{(k)})$ and $D$ is the Kullback-Leibler divergence (Cover and Thomas, 2006):

$$D(s^{(j)} \mid s^{(k)}) = \sum_{i=1}^{K} s_i^{(j)} \times \log_2 \frac{s_i^{(j)}}{s_i^{(k)}}$$

(2)

Novelty ($\mathcal{N}$) is defined as a document $s^{(j)}$'s reliable difference from past documents $s^{(j-1)}, s^{(j-2)}, \ldots, s^{(j-w)}$ in window $w$:

$$\mathcal{N}_w(j) = \frac{1}{w}\sum_{d=1}^{w} JSD(s^{(j)} \mid s^{(j-d)})$$

(3)

Resonance ($\mathcal{R}$) is defined as the degree to which future documents $s^{(j+1)}, s^{(j+2)}, \ldots, s^{(j+w)}$ conform to the Novelty of document $s^{(j)}$:

$$\mathcal{R}_w(j) = \mathcal{N}_w(j) - \mathcal{T}_w(j)$$

(4)

where $\mathcal{T}$ is the Transience of $s^{(j)}$:

$$\mathcal{T}_w(j) = \frac{1}{w}\sum_{d=1}^{w} JSD(s^{(j)} \mid s^{(j+d)})$$

(5)

Given the definitions outlined above, we can see that these information theoretic measures neatly translate into easily interpretable descriptions of the learning process over time. **Novelty** in our setup describes by how much the predictions of a given model at a particular training step differ from those which have come immediately before, indicating a substantial shift in model behavior. **Resonance**, on the other hand, considers to what extent this novelty persists in the system during subsequent training steps. This further allows the characterization of individual (per-experiment-level) signals as information dynamics profiles based on internal representation change. These information patterns can then be analyzed to see how the dynamics of the internals of a language model system evolve over time (i.e. during fine-tuning).

### 2.6 Signal processing

Due to the granularity of the experiments, the generated information signals are very long (as determined by batch size multiplied by number of training steps). As such, some processing must be done to analyze and interpret the signals meaningfully.

First, the first 160 and last 160 (i.e., the window size) documents are removed from the novelty and resonance signals. Second, following existing research into information dynamics (Nielbo et al., 2021a; Wevers et al., 2021; Nielbo et al., 2021b), non-linear adaptive filtering is performed to extract global trends in the novelty and resonance signals. In broad terms, the algorithm identifies a globally smooth trend signal by 'stitching' together locally best-fitting polynomials in overlapping partitions of the time series, allowing identification of broad trends while preserving local variations within the data. Following Riley et al. (2012), we define the span value (size of the partitions) by visually inspecting the results across a range of values to identify the best fit to extract the globally smooth trend across the different signals. In this study, this is done by comparing the smoothed signal produced by adaptive filters with varying span values to a moving average (see Appendix C.1 for an example). Based on this procedure, the span value for the partitions is set to 92.

## 3 Results

Figure 2 depicts the smoothed, normalized novelty and resonance signals for the 12 experiments in the fixed group (2a) and the optimized group (2b). Across both groups, the resonance signals show more frequent and periodic oscillations compared to the novelty signals. The trajectories of both novelty and resonance signals in the fixed group show a higher degree of similarity across experiments compared to those of the optimized group.

In the fixed group (Figure 2a), both novelty and resonance signals appear smoother and more coherent with slower oscillations, and we observe visible patterns that correlate across the different experiments. The novelty signals show closely aligned trajectories during the initial training phase, but begin to diverge after seeing approximately 20% of the documents. The divergence is apparent in the magnitude of the fluctuations, with some models having more or less pronounced variance. However, the overall direction of the changes - either increasing or decreasing - remains largely consistent across experiments. Though more variable from the outset, the resonance signals show similar patterns of divergence over time; they exhibit somewhat aligned trajectories in the initial training phase, but the magnitude of the oscillations grows more unsynchronized as the training progresses.

Figure 2: Normalized, smoothed novelty and resonance signals for experiments in (a) the fixed group and (b) the optimized group. The signals are visualized over fine-tuning time with the percentage of training documents seen by the model on the x-axis.

The situation is markedly different in the optimized group (Figure 2b), where the signals are more chaotic and noisy overall, with more rapid fluctuations and less apparent structure. Both novelty and resonance signals show high variability from the beginning of training and remain unsynchronized throughout. We observe less alignment across experiments, with more rapid fluctuations and no clear common direction of changes between experiments.

Figure 3 displays the novelty and resonance signals of the fixed group grouped by fine-tuning task. We can observe clear task-specific patterns in the trajectories of the signals, with high within-task alignment, especially for the MRPC task (Figure

3b). The same is not evident for the optimized group, nor do we find visible shared patterns in either group when grouping signals by model type (see Appendix B.1).

## 4 Discussion

Our findings reveal variations in information dynamics during the learning process across all experiments, suggesting that BERT models process and handle new information in distinct ways as they learn. Most notably, we observe a high degree of similarity in the signals from experiments in the fixed group. Despite divergences in magnitude, the overall directions of the changes in novelty and resonance remain largely consistent across exper-

Figure 3: Normalized, smoothed novelty (blue) and resonance (green) signals for experiments in the fixed group grouped by task. The highlighted lines depict the signals from fine-tuning on the (a) MNLI, (b) MRPC, and (c) SST-2 task, respectively. The transparent lines show the remaining signals, i.e., signals from those tasks not highlighted in each plot.

iments, suggesting somewhat stable information structure and shared underlying trends in the evolution of those signals over time. This contrasts with the optimized group, where the signals show more variability and noise across experiments, implying less consistency in information dynamics in that group. Contrary to previous methodologically related research in other data domains (e.g. Nielbo et al. (2021b) and Vrangbæk and Nielbo (2021)), this study does not find clear temporal change points in the information signals that correspond to key events, such as shifts in learning curves (see Appendix B.2).

The consistency we observe in the novelty signals in the fixed group suggests that, across model types, new information is being integrated in a stable and comparable way. The resonance curves show similar trends across models and tasks, indicating that when new information is introduced,

its influence tends to persist consistently across experiments. This illustrates a shared structure of learning dynamics, where the models steadily adapt to incoming training data in a similar manner. In contrast, while generally achieving better classification task performance (see Appendix B.2), the optimized group exhibits less consistent information integration. Frequent and high fluctuations in novelty signals in this group suggest that the models are encountering more abrupt changes in their internal representations, likely due to different optimal hyperparameters (e.g. learning rate or batch size). Resonance signals are also less uniform, implying that the influence of novel information on future representations is less predictable and more specific to the given experiment. These observations suggest that hyperparameter optimization introduces variability in how models process and retain information, possibly due to faster convergence, more aggressive adaptation, and divergent learning regimes across runs. However, it remains unclear whether these fluctuations reflect meaningful learning phenomena — such as adaptive capacity or sensitivity to task complexity — or are artifacts introduced by tuning. Distinguishing between the two remains a challenge and motivates future work involving finer-grained ablation studies and statistical analysis. Overall, these findings indicate that stability in training procedure (i.e., fixed hyperparameters) leads to more uniform information dynamics, while optimization increases variability in novelty and resonance, even if it may improve downstream task performance.

These results are aligned with prior work investigating fine-tuning dynamics in BERT models. For instance, as previously introduced, Hao et al. (2020) use divergence-based methods to assess shifts in attention patterns and find that fine-tuning affects the higher layers of BERT more substantially than lower layers. Their findings suggest that learning-induced changes tend to concentrate in specific architectural regions of the model and vary by downstream task — a conclusion that aligns with our observation that models under fixed training conditions exhibit consistent internal changes with observable task-specific patterns, while those under optimized regimes display greater variance. Given these earlier findings, the present study's focus on the prediction layer is a natural starting point for capturing salient representational changes during fine-tuning. However, while this level offers

tractable insight into the model's learning behavior, it may not fully capture the dynamics occurring in earlier layers. Extending the analysis to intermediate representations could provide a more nuanced understanding of how internal structures evolve across the network.

While the results may already conform with expert intuition about how models are learning over time, the explicitly information-theoretic approach can provide a new vocabulary and conceptual framework for explaining how and why certain learning dynamics occur during fine-tuning on different tasks.

For example, Figure 3 illustrates the information signals with fixed hyperparameters grouped by classification task. For all three tasks, there is an initial spike in novelty around 10% into training, indicating that significant, consecutive representational changes are occurring at this stage. This may reflect initial learning in the early stages of fine-tuning where the models make more sporadic or uninformed predictions, thus increasing novelty. Subsequently, novelty decreases, suggesting that the changes become more permanent, perhaps as the models have learned useful patterns from the training data. This is notably followed by a series of oscillations that manifest themselves consistently within each task, perhaps reflecting episodic shifts in representations as the models adjust to task-specific data.

The resonance signals show similarly pronounced regularity with structured, repeating resonance peaks, especially for the MRPC task (Figure 3b). This periodicity might emerge from uniform training dynamics across runs with fixed training regimes; the same types of examples tend to retain influence throughout training. The prominent resonance fluctuations in the MRPC task may correspond with overfitting tendencies observed in the learning curves of models fine-tuned on this task (see Appendix B.2). This suggests that certain training examples in MRPC repeatedly shape model behaviour, potentially leading to memorization rather than generalization.

These discussions highlight how the perspective introduced here offers not only exploratory or descriptive insights but also opens up for practical applications, such as change point detection. This may allow us to identify critical transitions in learning, e.g. sudden shifts in model behavior, convergence phases, or the onset of overfitting, poten-tially offering a more nuanced view of the training progress. While qualitative patterns suggest links between signal fluctuations and learning phenomena (e.g., spikes in novelty during early training), we do not currently quantify these relationships. The scope of this study is primarily descriptive and comparative; we focus on establishing the plausibility and interpretability of the proposed signals across training conditions. Future work could build on this foundation by investigating formal change point detection techniques or correlating signal dynamics with shifts in validation loss (Appendix B.2) to strengthen causal interpretations. We leave these directions for future research.

## 5 Conclusion

This paper presented a novel method demonstrating how information-theoretic signals can offer insights into the dynamics of how language models process and integrate information during fine-tuning. While traditional evaluation metrics provide static snapshots of model performance, our findings underscore the value of examining temporal learning dynamics to uncover how internal representations evolve over time. Across fixed training settings, models exhibit synchronous and structured changes, while optimized training regimes introduce greater variability, thus revealing how different learning conditions shape information flow.

For the purpose of this study, we focused only on BERT-style models, but the methods proposed here can be extended to other architectures. both the information-theoretic framework and the format of the GLUE benchmark can be model-agnostic, meaning that this analysis could feasibly be extended to different architectures, training regimes, and tasks. By quantifying how models react to and retain new information, this moves beyond performance outcomes to illuminate *how* models learn, not just how well. It captures the learning process as a sequence of representational shifts, offering a mathematical perspective on learning as continuous adaptation rather than discrete updates. Our work contributes a new layer of transparency to model behavior, bridging performance metrics with internal state changes, and advancing our understanding of learning as an unfolding, temporal process.

## Limitations

### Signal processing and window size

The generated information signals are inherently dependent on the chosen window size, as this defines the context for measuring 'surprise'. In this study, the choice of window size was intended to balance the trade-off between capturing sufficient context from the surrounding documents while maintaining computational feasibility. Though meaningfully defining an optimal window size for a problem as such remains a complex challenge, a sensitivity analysis (see Appendix C.2) showed that varying the window size within a small range had minimal impact on results. Still, all tested sizes were relatively short compared to the full signal. Future work could explore larger windows to examine long-term trends, though comparing distributions over broader spans may introduce limitations due to memory constraints in the current information-theoretic measures. Additionally, as previously discussed, the choice of adaptive filter span value was guided by visual inspection due to the lack of a standardized quantitative criterion for adaptive filter tuning. Though a range of values were tested for each experiment (Appendix C.1), its effects on signal smoothing could be explored more systematically in future work.

### Model and task diversity

The classification tasks were carefully selected to span a variety of differing scenarios. However, extending this work to include more complex classification problems, such as with imbalanced data or a wide number of classes, could offer additional insights. Likewise, our current work has been confined to English language tasks. While we found minimal differences between multilingual and monolingual BERT models, further investigation could clarify how language diversity shapes information dynamics. Similarly, while the models examined in this study have notable differences in architecture and training regimes, they all share the same BERT-style model at their core. Comparing information dynamics across more diverse model types could reveal alternative learning patterns and deepen our understanding of how different architectures integrate and retain information. While our model and task selection ensure a manageable comparison scope, extending this framework to other architectures (e.g., T5, GPT) and task types (e.g., generation, multilingual classification) would help assess the generalizability of information signals across broader learning paradigms.

## Ethics Statement

This study aimed to aid in opening the 'black box' of LLMs and enhance transparency by exploring the information dynamics in their internal representations. It takes an exploratory and analytical approach in nature and does not involve model deployment, private user data, or human subjects. The dataset used is publicly available and widely used in the research community. While our work contributes to model transparency research, it does not provide definitive explanations of model decisions. We caution against potential misuse, such as over-interpreting signals or applying our framework to justify opaque model behavior without sufficient validation. Finally, we must consider the environmental impact of our work, with 24 fine-tuning experiments and subsequent generation of information signals.

## Acknowledgments

## References

Alexander T. J. Barron, Jenny Huang, Rebecca L. Spang, and Simon DeDeo. 2018. Individuals, institutions, and innovation in the debates of the french revolution. *Proceedings of the National Academy of Sciences*, 115(18):4607–4612.

Tom B. Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, Sandhini Agarwal, Ariel Herbert-Voss, Gretchen Krueger, Tom Henighan, Rewon Child, Aditya Ramesh, Daniel M. Ziegler, Jeffrey Wu, Clemens Winter, Christopher Hesse, Mark Chen, Eric Sigler, Mateusz Litwin, Scott Gray, Benjamin Chess, Jack Clark, Christopher Berner, Sam McCandlish, Alec Radford, Ilya Sutskever, and Dario Amodei. 2020. Language models are few-shot learners.

Angelica Chen, Ravid Shwartz-Ziv, Kyunghyun Cho, Matthew L. Leavitt, and Naomi Saphra. 2025. Sudden drops in the loss: Syntax acquisition, phase transitions, and simplicity bias in mlms.

Andy Coenen, Emily Reif, Ann Yuan, Been Kim, Adam Pearce, Fernanda Viégas, and Martin Wattenberg. 2019. Visualizing and measuring the geometry of bert.

Thomas M Cover and Joy A Thomas. 2006. *Elements of Information Theory*, 2nd edition. John Wiley & Sons, Inc.

Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. 2019. BERT: Pre-training of deep bidirectional transformers for language understanding. In *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers)*, pages 4171–4186, Minneapolis, Minnesota. Association for Computational Linguistics.

William B. Dolan and Chris Brockett. 2005. Automatically constructing a corpus of sentential paraphrases. In *Proceedings of the Third International Workshop on Paraphrasing (IWP2005)*.

Ossama Embarak. 2023. Decoding the Black Box: A Comprehensive Review of Explainable Artificial Intelligence. In *2023 9th International Conference on Information Technology Trends (ITT)*, pages 108–113.

Yaru Hao, Li Dong, Furu Wei, and Ke Xu. 2020. Investigating learning dynamics of bert fine-tuning. In *AACL*.

Vikas Hassija, Vinay Chamola, Atmesh Mahapatra, Abhinandan Singal, Divyansh Goel, Kaizhu Huang, Simone Scardapane, Indro Spinelli, Mufti Mahmud, and Amir Hussain. 2024. Interpreting Black-Box Models: A Review on Explainable Artificial Intelligence. *Cognitive Computation*, 16(1):45–74.

Evan Hernandez and Jacob Andreas. 2021. The low-dimensional linear geometry of contextualized word representations. In *Proceedings of the 25th Conference on Computational Natural Language Learning*, pages 82–93, Online. Association for Computational Linguistics.

Ross Deans Krisensen-McLachlan, Rebecca M.M. Hicke, Márton Kardos, and Thunø Mette. 2024. Context is Key(NMF). *CHR 2024: Computational Humanities Research Conference*, pages 829–847.

Nelson F. Liu, Matt Gardner, Yonatan Belinkov, Matthew E. Peters, and Noah A. Smith. 2019. Linguistic knowledge and transferability of contextual representations. In *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers)*, pages 1073–1094, Minneapolis, Minnesota. Association for Computational Linguistics.

David J. C. MacKay. 2002. *Information Theory, Inference & Learning Algorithms*. Cambridge University Press, Cambridge.

Amil Merchant, Elahe Rahimtoroghi, Ellie Pavlick, and Ian Tenney. 2020. What happens to BERT embeddings during fine-tuning? In *Proceedings of the Third BlackboxNLP Workshop on Analyzing and Interpreting Neural Networks for NLP*, pages 33–44, Online. Association for Computational Linguistics.

Kristoffer L. Nielbo, Rebekah Brita Baglini, Peter Bjerregaard Vahlstrup, Kenneth C. Enevoldsen, Anja Bechmann, and Andreas Roepstorff. 2021a. News information decoupling: An information signature of catastrophes in legacy news media. *CoRR*, abs/2101.02956.

Kristoffer L. Nielbo, Frida Hæstrup, Kenneth C. Enevoldsen, Peter B. Vahlstrup, Rebekah B. Baglini, and Andreas Roepstorff. 2021b. When no news is bad news – Detection of negative events from news media content. *arXiv:2102.06505 [cs]*. ArXiv: 2102.06505.

Jason Phang, Haokun Liu, and Samuel R. Bowman. 2021. Fine-Tuned Transformers Show Clusters of Similar Representations Across Layers. ArXiv:2109.08406 [cs].

Alec Radford, Jeffrey Wu, Rewon Child, David Luan, Dario Amodei, and Ilya Sutskever. 2019. Language models are unsupervised multitask learners. *OpenAI blog*, 1(8):9.

Michael A. Riley, Scott Bonnette, Nikita Kuznetsov, Sebastian Wallot, and Jianbo Gao. 2012. A tutorial introduction to adaptive fractal analysis. *Frontiers in Physiology*, 3.

Baptiste Rozière, Jonas Gehring, Fabian Gloeckle, Sten Sootla, Itai Gat, Xiaoqing Ellen Tan, Yossi Adi, Jingyu Liu, Romain Sauvestre, Tal Remez, Jérémy Rapin, Artyom Kozhevnikov, Ivan Evtimov, Joanna Bitton, Manish Bhatt, Cristian Canton Ferrer, Aaron Grattafiori, Wenhan Xiong, Alexandre Défossez, Jade Copet, Faisal Azhar, Hugo Touvron, Louis Martin, Nicolas Usunier, Thomas Scialom, and Gabriel Synnaeve. 2024. Code llama: Open foundation models for code.

Victor Sanh, Lysandre Debut, Julien Chaumond, and Thomas Wolf. 2020. Distilbert, a distilled version of bert: smaller, faster, cheaper and lighter.

Julien Simon. 2021. Large Language Models: A New Moore's Law?

Richard Socher, Alex Perelygin, Jean Wu, Jason Chuang, Christopher D. Manning, Andrew Ng, and Christopher Potts. 2013. Recursive deep models for semantic compositionality over a sentiment treebank. In *Proceedings of the 2013 Conference on Empirical Methods in Natural Language Processing*, pages 1631–1642, Seattle, Washington, USA. Association for Computational Linguistics.

Ian Tenney, Patrick Xia, Berlin Chen, Alex Wang, Adam Poliak, R. Thomas McCoy, Najoung Kim, Benjamin Van Durme, Samuel R. Bowman, Dipanjan

Das, and Ellie Pavlick. 2018. What do you learn from context? Probing for sentence structure in contextualized word representations.

Elena Voita and Ivan Titov. 2020. Information-Theoretic Probing with Minimum Description Length. ArXiv:2003.12298 [cs].

Eva Elisabeth Houth Vrangbæk and Kristoffer Laigaard Nielbo. 2021. *Composition and Change in De Ciuitate Dei: A Case Study of Computationally Assisted Methods*, volume 14: Augustine of Hippo's De ciuitate Dei: Content, Transmission, and Interpretations of *Studia Patristica*, pages 149–164. Peeters.

Alex Wang, Yada Pruksachatkun, Nikita Nangia, Amanpreet Singh, Julian Michael, Felix Hill, Omer Levy, and Samuel R. Bowman. 2020. Superglue: A stickier benchmark for general-purpose language understanding systems.

Alex Wang, Amanpreet Singh, Julian Michael, Felix Hill, Omer Levy, and Samuel Bowman. 2018. GLUE: A multi-task benchmark and analysis platform for natural language understanding. In *Proceedings of the 2018 EMNLP Workshop BlackboxNLP: Analyzing and Interpreting Neural Networks for NLP*, pages 353–355, Brussels, Belgium. Association for Computational Linguistics.

Melvin Wevers, Jan Kostkan, and Kristoffer L. Nielbo. 2021. Event flow – how events shaped the flow of the news, 1950-1995.

Adina Williams, Nikita Nangia, and Samuel Bowman. 2018. A broad-coverage challenge corpus for sentence understanding through inference. In *Proceedings of the 2018 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long Papers)*, pages 1112–1122, New Orleans, Louisiana. Association for Computational Linguistics.

Thomas Wolf, Lysandre Debut, Victor Sanh, Julien Chaumond, Clement Delangue, Anthony Moi, Pierric Cistac, Tim Rault, Remi Louf, Morgan Funtowicz, Joe Davison, Sam Shleifer, Patrick von Platen, Clara Ma, Yacine Jernite, Julien Plu, Canwen Xu, Teven Le Scao, Sylvain Gugger, Mariama Drame, Quentin Lhoest, and Alexander Rush. 2020. Transformers: State-of-the-art natural language processing. In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing: System Demonstrations*, pages 38–45, Online. Association for Computational Linguistics.

Yichu Zhou and Vivek Srikumar. 2022. A Closer Look at How Fine-tuning Changes BERT. In *Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 1046–1061, Dublin, Ireland. Association for Computational Linguistics.

Liu Zhuang, Lin Wayne, Shi Ya, and Zhao Jun. 2021. A robustly optimized BERT pre-training approach with post-training. In *Proceedings of the 20th Chinese National Conference on Computational Linguistics*, pages 1218–1227, Huhhot, China. Chinese Information Processing Society of China.

# Appendices

## A Additional methods

### A.1 Hyperparameters for the optimized group

Before fine-tuning for experiments in the optimized group, we performed hyperparameter tuning, as described in the paper. The resulting hyperparameter configurations can be found in Table 2. Hyperparameters not specified in the table were kept at default values.

| Model | Task | $N$ epochs | LR | Batch size |
|---|---|---|---|---|
| BERT | MNLI | 3 | $3e^{-5}$ | 16 |
| BERT | MRPC | 2 | $5e^{-5}$ | 32 |
| BERT | SST-2 | 2 | $2e^{-5}$ | 16 |
| distilBERT | MNLI | 4 | $5e^{-5}$ | 32 |
| distilBERT | MRPC | 3 | $3e^{-5}$ | 16 |
| distilBERT | SST-2 | 2 | $5e^{-5}$ | 64 |
| roBERTa | MNLI | 4 | $2e^{-5}$ | 32 |
| roBERTa | MRPC | 4 | $2e^{-5}$ | 64 |
| roBERTa | MRPC | 4 | $5e^{-5}$ | 64 |
| mBERT | MNLI | 4 | $5e^{-5}$ | 32 |
| mBERT | MRPC | 3 | $5e^{-5}$ | 64 |
| mBERT | SST-2 | 3 | $2e^{-5}$ | 64 |

Table 2: Hyperparameter configurations for each experiment in the optimized group. LR is the learning rate.

## A.2 Model architectures and pre-training details

In Table 3, we highlight some of the main differences between the four models in terms of architecture and pre-training details.

| Model | $N$ layers | $N$ parameters | $N$ languages |
|---|---|---|---|
| BERT | 12 | 110M | 1 |
| distilBERT | 6 | 66M | 1 |
| roBERTa | 12 | 125M | 1 |
| mBERT | 12 | 110M | 102 |

Table 3: Overview of architecture and training details for pre-trained versions of BERT, distilBERT, roBERTa, and mBERT.

## B Additonal results

### B.1 Grouped signals

To explore patterns in the extracted information signals, different groupings of the signals were visualized. As discussed in the paper, the analysis revealed task-specific patterns in the information signals from the experiments in the fixed group. In Figure 4, the information signals from the

optimized group are shown grouped by task. All subfigures display all the same signals; however, each subfigure highlights the novelty and resonance signals for a respective task, while the remaining signals are depicted in transparent lines for comparison.



Figure 4: Normalized, smoothed novelty (blue) and resonance (green) signals for experiments in the optimized group grouped by task. The highlighted lines depict the signals from fine-tuning on the (a) MNLI, (b) MRPC, and (c) SST-2 task, respectively. The transparent lines show the remaining signals, i.e., signals from those tasks not highlighted in each plot.

Similarly, the information signals from the experiments were grouped by model type to investigate potential patterns. This is depicted in Figure 5, with each row of subfigures highlighting the signals of the four different models, respectively. The left column shows experiments from the fixed group, and the right column shows experiments from the optimized group.

Figure 5: Normalized, smoothed novelty (blue) and resonance (green) signals grouped by model. The left column shows the fixed group and the right column shows the optimized group. Each row corresponds to a model: (a–b) BERT, (c–d) distilBERT, (e–f) roBERTa, and (g–h) mBERT. Highlighted lines show signals for each model-task combination; transparent lines show the rest.

## B.2 Learning curves

In Figure 6, the learning curves for the various experiments are presented, illustrating the models' performances on the classification tasks during the fine-tuning process. Each subfigure represents an experiment, displaying the learning curves for each of the models fine-tuned on a task. The purple line represents validation accuracy, the red line represents validation loss, and the yellow line represents training loss. Note that differing training durations in the optimized group led to uneven checkpoint sampling across experiments. As a consequence, some plots — such as those for roBERTa — are missing or incomplete (e.g., if training terminated before enough checkpoints were saved).

## C  Sensitivity analyses

### C.1  Defining the adaptive filter span

As discussed in Section 2.6, we follow the proposed method for defining the span value for the adaptive filter (Riley et al., 2012); namely, visual inspection of the fit of the smoothed signal produced by varying span values. Figure 7 displays an example of this.

### C.2  Defining the window size

As mentioned in the Limitations, a sensitivity analysis was also performed to investigate the effect of varying the window size in which to calculate the information signals. An example of this can be seen in Figure 8.

(a)



(b)



Figure 6: Learning curves for experiments in (a) the fixed group and (b) the optimized group. The red line represents the validation loss, the yellow line represents the training loss, and the purple line represents the validation accuracy.

227

Figure 7: Example of the effect of different span values for the adaptive filter. The signal depicted here is the novelty signal from BERT fine-tuned on the MRPC task with fixed hyperparameters. The grey line depicts the original, unsmoothed novelty signal. The light blue line depicts the novelty signal's moving average (w=10000). The dark blue line depicts the smoothed signal from the adaptive filter using span values of (a) 32, (b) 56, and (c) 128, respectively. All signals are normalized.



Figure 8: Example of the effect of varying the window size for which to calculate novelty, transience, and resonance in. The signal depicted is the normalized, smoothed novelty signal from distilBERT fine-tuned on the MNLI task with fixed hyperparameters. The different lines represent different window sizes (80, 160, 320)

228

# GUIDE: A Framework for Improving Functional Software Test Descriptions with Language Models

**Mathis Ronzon**
Alten Labs, Rennes, France
mathis.ronzon2@alten.com

**Thierry Roger**
Alten Labs, Rennes, France
thierry.roger@alten.com

**Zoltan Miklos**
Univ Rennes CNRS IRISA
Rennes, France
zoltan.miklos@irisa.fr

**Annie Foret**
Univ Rennes CNRS IRISA
Rennes, France
annie.foret@irisa.fr

## Abstract

Functional software testing is essential to ensure that software meets user expectations. Our ambition is to enable business experts who have extensive domain knowledge but limited software engineering competences, to realize the functional software tests, through formulating test case descriptions in natural language. To meet this challenge, we propose a framework called GUIDE (Guided User-driven Interactive Description Enhancement), which leverages language models to improve functional software test descriptions written in natural language. Our framework implements an intermediate step based on a structured language (Gherkin[1]) that is a language widely used for software tests. We translate test descriptions written by the business expert to this language using language models. We automatically evaluate the quality of the description based on the generated Gherkin. When this quality is insufficient, GUIDE initiates an interactive and personalized assistance process, delivering targeted advice to help business experts enhance and improve their test case descriptions. We evaluated our approach through a case study based on test cases for a human resources management related software, written in French. We recorded a 26% decrease in the average number of descriptions required per test objective to reach the desired quality level thanks to the advice generated.

## 1 Introduction

Software testing plays a crucial role in the quality and longevity of IT applications. However, a persistent divide between developers and end-users often complicates this essential task. While developers master the code and technical specifics, the business expert is the one who know the functional requirements. Nevertheless, they do not always have the tools or the language to express their expectations that is precisely understandable by the developers. This dissonance can compromise test reliability and software quality.

The advent of language models, capable of generating code from natural language instructions, is a promising solution for test production (Tufano et al., 2021; Xie et al., 2023). However, this approach is still mainly accessible to people with solid programming expertise, thus excluding many business experts. The latter also encounter difficulties in interacting with language models: they often give up too early when faced with a lack of understanding of the model, or formulate erroneous expectations based on dynamics specific to human interaction (Zamfrescu-Pereira et al., 2023). Moreover, each individual writes natural language text differently, even when the objective is identical (Weigelt et al., 2020). In parallel, some recent benchmarks (Jimenez et al., 2024) now include expert-verified and human-annotated versions of problem descriptions, acknowledging that instructions written by non-expert users are often insufficient to fully capture the functional intent behind a request.

In light of these challenges, it becomes essential to explore novel approaches that qualify user input and provide actionable guidance to improve it before giving it to a code generation pipeline. To this end, we propose a novel framework that we call GUIDE (*Guided User-driven Interactive Description Enhancement*), that aims to guide business experts in improving the quality of functional software test descriptions. The main objective of GUIDE is to enhance the clarity and precision of natural language test descriptions. More specifically, these descriptions are produced by a business expert without any specific rules being imposed,

---

[1] https://cucumber.io/docs/gherkin/reference/

Figure 1: Overview of the GUIDE framework. In blue are the parts that the user is asked to perform, and in black those that are automated. Firstly, the user produces a description of a functional test, which is then translated into the Gherkin language using a language model. The quality of the test description is then assessed using this code. If the quality is judged to be insufficient, advice is generated based on the test description and the Gherkin code, enabling the user to modify the description himself.

and seek to explain how the system works according to a particular functionality. Our approach is based on three key components (Figure 1): (1) a quality criterion leveraging the automatic translation of descriptions into an intermediate representation language (Gherkin), (2) a classification system to evaluate the semantic similarity between the description and its translation, and (3) an interactive advice generation mechanism to guide users in refining their descriptions. A complete example of such a test description, its translation into Gherkin and their similarity according to our scale can be found in Table 1.

To evaluate the effectiveness and relevance of GUIDE, we seek to address the following research questions:

- **RQ1** : To what extends automatic translation of a test description into an intermediate language can constitute an appropriate quality criterion?

- **RQ2** : Can a small language model manage to understand the semantic similarities between a test description and its translation into an intermediate language?

- **RQ3** : Does an interactive process using automatic advice generation can help a user improve the quality of his description?

Thanks to the participation of 60 people, spread

over three labelling campaigns and two production campaigns, we have been able to evaluate the effectiveness of GUIDE. The results show that 70% of the advice generated is considered relevant by users. In addition, we observed a measurable improvement in their ability to comply with the quality criteria, with a 26% reduction in the average number of descriptions needed per test objective to achieve the required level of quality.

The rest of the paper is organized as follows. We start by discussing the related works in Section 2. In Section 3, we introduce the three key components of our GUIDE framework in detail. Then, we discuss the practical implementation of GUIDE in a real-world experiment, where business experts interact with software under test conditions, that we describe in Section 4. We conclude the paper in Section 5.

## 2 Related Work

We review contributions related to large language models for code generation, techniques for refining ambiguous or incomplete user input, and methods for assessing the quality of natural language descriptions through text classification. We conclude by positioning our approach, GUIDE, in relation to these works.

### 2.1 LLM-based Code Generation

Large language models (LLMs) have emerged as a powerful tool for various code-related tasks, including program synthesis (Austin et al., 2021), bug fixing (Zubair et al., 2024) or program testing (Xiong et al., 2023). Through extensive pre-training, they recognize patterns, comprehend context, and generate coherent and contextually relevant code snippets.

In software testing, the use of natural language as an entry point remains limited. Approaches such as AthenaTest (Tufano et al., 2021) or A3Test (Alagarsamy et al., 2024) rely mainly on source code to generate tests, while ChatUniTest (Xie et al., 2023) uses a prompt composed mainly of code fragment.

### 2.2 Refining user input

The task of asking users to reformulate or modify their output is receiving increasing attention in the fields of information retrieval and dialogue systems. For example, Wang and Li (2021) propose a method based on question templates to help users clarify their requests. They use a genera-

Table 1: An example of a functional test description written by a non AI expert, translated into Gherkin code by an LLM and its similarity given by a human annotator using our similarity scale (Appendix A) (Example from our use case. The original description was written in French and can be found in the first row of Table 6)

| Description | Gherkin | Similarity Label |
|---|---|---|
| Click on 'OPEN'. Check that the file selection window opens. Select a file with the '.csv' extension and check that the search appears in the software. | **Given** the software is open<br>**When** I click on "OPEN"<br>**Then** the file selection window opens<br>**And** I select a file with the extension ".csv"<br>**And** the search appears in the software | VERSIM |

tion model based on a Transformer. Eberhart and McMillan (2022) propose a new method that uses a task extraction algorithm to identify aspects of the query and follows a rule-based procedure to generate questions.

In code generation, dealing with ambiguous user requirements has received more attention. Some pipeline such as QualityFlow (Hu et al., 2025) have integrated the evaluation of the quality of natural language requests and offers self-improvement mechanisms to reformulate instructions without the aid of the user. Other methods such as ClarifyGPT (Mu et al., 2023) or CodeClarQA (Li et al., 2022) questions the user to clarify ambiguities. In the case of ClarifyGPT, the ambiguity of a requirement is detected when several generations of code produced from the same instruction lead to different behaviours for the same input. As for CodeClarQA, it always asks questions, but has no system for assessing this ambiguity.

### 2.3 Text classification

To address the issue of ambiguous or incomplete test descriptions, recent research has focused on automated classification of textual quality. Traditional methods rely on logical rule-based systems that detect key phrases or syntactic patterns indicative of test completeness (Ormandjieva et al., 2007). Although effective in well-structured scenarios, these systems lack the flexibility to handle the linguistic diversity present in natural language inputs. Another method involves supervised learning based on standard metrics, representing the criteria that an expert takes into consideration when assessing the quality of requirements (Parra et al., 2015).

### 2.4 GUIDE positioning

We have seen that in current approaches to code generation, consideration of the quality of user input often comes after a long and costly process, once the initial generation has already been carried out. In contrast, GUIDE seeks to intervene upstream, validating the quality of the request as soon as it is created, to ensure that all the elements required for correct generation are present.

In addition, when ambiguity is detected, some approaches ask questions, sometimes of a technical nature, to fill in the missing information. GUIDE adopts a different strategy: it allows the user to modify the description directly, making implicit information explicit. This process is based on the generation of targeted advice, offered to the user in a non-binding way. In this way, the user retains control of their production, while being guided to improve it progressively.

## 3 The framework GUIDE

In this section, we will present our framework. In order to do so, we start by discussing quality criterion that we have chosen to discard the description. Then, we will look at the method used to assess automatically the quality of a test description using its translation into Gherkin code. Finally, we will present the method used to guide the user in the process of improving his production.

### 3.1 Quality criterion

Our approach is based on the use of a quality criterion to filter test descriptions according to their relevance. Our aim is to develop a method that retains only those test descriptions that contain all the elements necessary for a language model to both generate and verify the corresponding test. Thus,

the quality criterion aims to evaluate the ability of the language model to restore all the information contained inside the description.

Thus, to assess the language model's ability to understand the description, we use it in a simple task: information reorganisation. Specifically, we ask the language model to structure the information present in the test description. To do so, we use a intermediate language that meets this requirement perfectly: Gherkin. Gherkin is a human-understandable specification language used in the development method known as BDD (Behavior-Driven Development).

We will therefore base our quality criterion on the similarity level between the test description and its translation into Gherkin code by a language model. More specifically, it is based on a similarity scale detailed in Appendix A, which is composed of five labels:

- Similar : COMPSIM, VERSIM

- Different : SOMSIM, VERDIFF, COMPDIFF

By aggregating these labels into two distinct groups, we define our quality criterion: a description is considered of good quality if it obtains a similarity label of COMPSIM or VERSIM. Event though our quality criterion does not take advantage of all the nuance offered by the five similarity labels, we have choose to retain this scale. This additional granularity proves invaluable during manual labelling campaign, as it enables human annotators to better express their perception of the similarity. This finer distinction encourages a more precise and nuanced assessment.

## 3.2 Supervised Learning Quality Assessment

Two main approaches can be used to classify a text according to a given label: logical methods, based on rules or feature extraction, and approaches based on language models. However, in our case, the test descriptions present a high degree of lexical and structural heterogeneity, in addition to a dense technical content. These characteristics make logical methods ineffective, as they are too rigid and not very adaptable to the variability of the data.

We therefore opted to use language models, which are better able to capture the subtleties of natural language, even in a technical context. Given that our Framework is likely to handle sensitive or confidential data, we have deliberately restricted our choice to compact models that can be run locally, without depending on remote services.

## 3.3 Interactive Improvement Process

GUIDE not only qualifies the description written by the user, but also provides guidance when it is not of good quality. This guidance is intended to suggest possible changes that the user could make. We have chosen to use this form because we want the user to have a choice of modification throughout the procedure.

To produce it, we will use the description and its associated Gherkin code and give it, using a pre-defined prompt, to a language model. In order not to introduce our own bias into this generation of advice, and the evaluation of the prompts used to generate them, we have decided to use an automatic optimisation methods to find the best prompt, in particular Beam Search (Pryzant et al., 2023). This method is based on a starting prompt, a scoring metric and a method for generating several variations of a prompt to explore the space of available prompts.

The basic prompt is structured in three main parts: two dedicated slots for inserting the test description and its translation into Gherkin; another slot used to insert context of the software under test to maximise the relevance of the advice generated; an explicit sentence tells the language model the expected objective as well as the constraints to be respected (clarity, consistency and respect for the context).

For the mutation prompt, we ask the language model to produce three variations of a prompt, while retaining the meaning and the three slots reserved for description, Gherkin and context. These variations aim to explore different formulations while maintaining the structure of the task. An example of a basic prompt and a mutation prompt can be found at Appendix C.

The score metric is based on a realistic approach, aimed to simulate the behaviour of a user using the advice generated by a prompt. Using a dataset of test descriptions classified as of bad quality, and the prompt submitted for evaluation, we will produce advice for each of them. This advice, combined with the initial description and the context of the software under test, is then provided to a language model which simulates a user by producing a new version of the description, incorporating the suggested recommendations. This new descriptions,

together with its translation into Gherkin, is submitted to the classifier for quality evaluation. We calculate the prompt score by measuring the percentage of descriptions that, after modifications, pass the quality criterion.

# 4 Software Testing by business expert : Case study

In order to evaluate GUIDE in a real-life application, we chose to involve users, not necessarily with a background in IT, in the task of writing test descriptions. To this end, we set up three complementary campaigns as illustrated in Figure 2:

- A labelling campaign which has the objective to evaluate the relevance of Gherkin productions automatically generated by a language model.

- A production campaign which aims to observe how users behave when writing descriptions, measuring in particular their ability to produce content in line with our quality criterion.

- An advice campaign that seeks to assess user satisfaction with the use of advice produced by our method.

For the sake of simplicity, we decided to ask for the descriptions to be produced in French for all the campaigns. This allowed us to have more people available to take part in the campaigns.

## 4.1 Software under test

The software to be tested is called "Esco Explorer". It is a tool for displaying a graph in the form of an Acyclic Guided Graph (AGG) based on the occupations/skills given by the Esco ontology (Appendix B). Esco is a European classification of skills, competencies, qualifications and professions. The system identifies and categorizes skills, competences, qualifications and occupations relevant to the EU labor market, education and training, in 25 European languages. The system provides occupational profiles showing the relationships between occupations, skills, competences and qualifications; it functions like a dictionary.

During software development, a test plan consisting of 68 tests divided into 9 categories was produced. We used this same test breakdown for the rest of this section, enabling us to indicate a category and a precise test goal to the user, to help

them write their description. An example of a category and its associated test purpose is: Category – Node Information, with the Test purpose – Display an optional job for a skill.

## 4.2 Campaigns

**Labelling Campaign**

Prior to the various experiments, five users were asked to write one natural language description per test purpose, resulting in five complete test plans. Based on these plans, we generated a manual labeling campaign aimed at assessing the quality of the descriptions via their correspondence with an automatically generated Gherkin code.

In concrete terms, each annotator was assigned a test plan, in which he or she had to select a description, consult the corresponding Gherkin code, generated by a language model, and then evaluate the similarity between the two elements. This evaluation was carried out using our similarity scale.

**Production Campaign**

During a second campaign, users will have to write test description for each test purpose themselves. For each test purpose, users will have to produce a test description then label the similarity between the description and its translation in Gherkin. If their description doesn't respect our quality criterion, they will have to modify their description and redo the labeling process.

In order to avoid blocking users when faced with cases they consider too complex, we have left open the possibility of changing the test to be described even when the quality criterion has not been met. However, to guarantee a minimum of reformulation effort, each user was required to propose at least two attempts to improve his initial description before having the possibility to abandon and do a new test.

**Advice Campaign**

Using the same protocol as the production campaign, this time we decided to add the tips generated from the prompt designed in Section 3.3. More specifically, when the user indicates that their description is not of sufficient quality, we offer them the advice generated from their description and the Gherkin. The user can then take these tips into account, or not use them if they don't find them interesting. This is indicated by two labels.

**Gherkin Generation**

In the interests of data governance and in order to guarantee local execution without dependency on

Figure 2: Overview of the different campaigns. The blue color represents the labelling campaign, where users assess the similarity between a previously written description and its Gherkin translation. The orange indicates the production campaign, in which users write a description based on a given test purpose, evaluate its quality, and revise it if necessary. Finally, the green corresponds to the advice campaign, where users receive guidance generated from their initial description and its Gherkin translation to help them improve their text.

external services, we opted to use the Mistral V0.3-7b language model (Jiang et al., 2023), quantised in 4 bits. The prompt used to generate Gherkin was manually optimised using a set of descriptions considered to be of good quality.

### 4.3 GUIDE Implementation

#### 4.3.1 Quality Classifier

In order to evaluate several possible classifiers based on language models, we have chosen to use CamemBERT v2 (Antoun et al., 2024) and SomlLM2-135M (Allal et al., 2025) as a basis. CamemBERT v2 is a robust and high-performance reference for automatic language processing tasks in French. SmolLM2-135M is a more recent, lightweight model, that is recognised for its good general capabilities despite its small size. This choice makes it possible to combine confidentiality, linguistic performance and operational efficiency.

Since SmolLM2 is trained exclusively on English data, we explored its viability by automatically translating our dataset into English using the opus-mt-fr-en (Tiedemann et al., 2023; Tiedemann and Thottingal, 2020) model. This model allows efficient conversion of descriptions and Gherkins written in French into English, ensuring consistent basis for training.

To train our classification models, we used the

data collected from two campaigns described in Section 4.2: the labelling campaign and the production campaign. This resulted in a dataset of 1522 test description / Gherkin code pairs. Each pair was transformed into a single classifier input using explicit keywords (Description:, Gherkin:) and a [SEP] separator token.

Based on our quality criterion, we converted the original similarity labels into binary quality labels: 885 instances were labeled as SIMILAR (indicating sufficient quality), and 637 as DIFFERENT. The dataset was split into a 70/30 ratio for training and validation. All models were trained with 3 epochs and a batch size of 8.

To complement this training set, we also constructed a separate test set of 68 description/Gherkin pairs. These descriptions were written by a single user not involved in the previous datasets, and each pair was annotated five times. The final label for each instance was computed as an average label, using the method described in Section 4.4.1.

Table 2: Performance of classification models in training and evaluation

| Model | Training Accuracy | Evaluation Accuracy | Test Accuracy |
|---|---|---|---|
| CBERT-Mix | 81.78 | 74.55 | 60.29 |
| Smol-FR | 98.97 | 73.52 | 70.59 |
| Smol-EN | 98.97 | 72.23 | 58.82 |

The training results for the different models can be found inside Table 2. The models based on SmolLM2 show better accuracy during training, but this superiority is not reflected on the evaluation set, where the performance is similar to that of the models based on CamemBERT. Furthermore, translating the data into English did not bring any significant improvement in terms of accuracy.

Of the models evaluated, only one managed to maintain good accuracy over the test set: Smol-FR. The other two models showed a significant drop in performance, indicating a more limited ability to generalise. As a result, we have chosen Smol-FR for the rest of our experiments.

#### 4.3.2 Advice Prompt

Using this classifier, we were able to launch the BeamSearch algorithm to produce the prompt used to generate the advice. We retrieved the first paragraph of Section 4.1 to be used as the context of the application under test inside the advice prompt.

As for the description dataset required for the score metric, we selected the descriptions from the production campaign (presented in the in Section 4.2), which were identified as being of insufficient quality by the classifier as well as the user who produced the description.

In total, our method evaluated 45 prompts each of which was evaluated with 218 data items. The selected prompt (found in Listing 3) received a score of 0.79.

## 4.4 Data Analysis

In this section, we will analyse the data from the three campaigns presented above. We will seek to answer the various research questions we had, as well as assessing the usefulness of GUIDE.

### 4.4.1 Gherkin Generation ability

**Mean Label**

A total of 21 people took part in the labelling campaign. Of the five test plans proposed, three were annotated by five people, while the other two were annotated by three people. Since several annotators evaluated the same pairs of data (test description and Gherkin code generated by an LLM), it is necessary to assign a consensus similarity label of each piece of data. To achieve this, we adopt a majority voting approach.

Using the decomposition of the five similarity label according to our quality criterion, we will look at the group with the most labels. Inside this majority group, if one label stand out with a clear majority, it is selected. Otherwise, we proceed to average the labels to select the most representative.

**Model capability**

Looking at the distribution of average labels obtained for each test plan (Table 3), we can see that some users, with no prior knowledge of the Gherkin language or the quality criterion used, manage to produce descriptions that directly satisfy this criterion. However, this success is not homogeneous: other test plans present initial descriptions whose quality is insufficient according to our quality criterion.

Despite these disparities, one encouraging point stands out: no test plan is completely misunderstood by the language model. This illustrates the robust ability of the selected model (MistralV0.3-7b) to interpret even descriptions from non-expert authors, and to produce usable Gherkin translations.

Table 3: Distribution of the similarity label depending of the Test plan.

| Test Plan | COMPSIM | VERSIM | SOMSIM | VERDIFF | COMPDIFF |
|---|---|---|---|---|---|
| **1** | 12 | 11 | 15 | 11 | 19 |
| **2** | 43 | 16 | 5 | 3 | 1 |
| **3** | 4 | 18 | 22 | 19 | 5 |
| **4** | 7 | 28 | 15 | 12 | 6 |
| **5** | 26 | 37 | 2 | 3 | — |

That said, a qualitative analysis of the comments left by annotators allows us to distinguish two main sources of error in Gherkin generation: problems of structure, linked to poor syntactic or logical organization of the generated code, and problems of ambiguity, due to an incomplete or poorly formulated initial description. This ambiguity is due to the annotator, who did not necessarily understand the test description correctly, as he pointed out in his commentary. It is therefore a semantic ambiguity, linked to imprecise wording or wording that is open to several interpretations in the initial description. Table 4 shows, for each test plan, the proportion of descriptions identified as having these two types of problem.

Table 4: Distribution of Structure and Ambiguity Errors in Low-Quality Descriptions per Test Plan (as Labeled by Annotators)

| Test Plan | Structure (%) | Ambiguity (%) |
|---|---|---|
| **1** | 90.2 | 58.6 |
| **2** | 76.3 | 88.9 |
| **3** | 89.2 | 73.0 |
| **4** | 90.7 | 65.8 |
| **5** | 55.2 | 51.8 |

We find that, in the majority of cases, failures to meet our quality criterion stem first and foremost from problems with the model's structuring of the Gherkin, with rates exceeding 90% in some shots. However, these structural errors are often exacerbated by poorly constructed initial descriptions, as shown by the high rate of ambiguity problems - reaching 88.9% in plan 9. This twofold observation highlights both the current limitations of the language model in correctly handling Gherkin's syntactic constraints, and the need to support users in improving the clarity and completeness of their descriptions.

#### 4.4.2 User behaviour

Based on the results presented in Table 8, we observe that all users needed to revise their descriptions at least once, and in many cases several times, before reaching a level that matched our defined quality criterion. This reinforces the idea that generating a high-quality test description is not straightforward, especially for non-expert users.

However, the low number of abandons suggests that the effort required to improve a description is not perceived as excessive. In particular, only 36 abandons were recorded out of 680 attempts, indicating that most users were willing to iterate to reach the expected quality level.

Each user had to produce descriptions for 9 different test categories, in a fixed order that was identical for everyone. This enabled us to observe a potential progression. However, no clear trend emerged. We find this to be the case even when we split each categories into two equal halves (Table 7). No systematic improvement dynamic can be observed. So we don't need to take into account a history for each user in our GUIDE framework.

#### 4.4.3 Advice capability

In total, 125 advice have been generated, with an overall satisfaction rate of 70%. This result indicated that the majority of users were positive about the usefulness of the advice provided.

The analysis also shows a reduction in the average number of descriptions per test: this drops from 1.75 (observed during the initial production campaign) to 1.29 during this campaign (Table 9). This reduction suggests that the advice makes it easier to achieve the quality criteria, thereby reducing the number of iterations required.

In addition, result in Table 5 reveal a marked difference between advice that is considered relevant and advice considered uninteresting. More specifically, advice perceived as useful is significantly more associated with improvements in the quality of the description. This trend suggests that the perceived quality of the advice has a direct influence on the user's ability to refine their description, thereby reinforcing the effectiveness of GUIDE's interactive process.

## 5 Conclusion

This work introduced GUIDE (Guided User-driven Interactive Description Enhancement), a framework that improves the quality of test descriptions

Table 5: Improvement of the similarity label depending on whether the advice has been deemed relevant by the user.

| Improvement | Interesting | Not Interesting |
|---|---|---|
| Upgrade | 51 | 11 |
| Constant | 34 | 24 |
| Downgrade | 2 | 3 |

written by business experts through an interactive process. By leveraging Gherkin as an intermediate representation, GUIDE effectively assesses description quality and provides personalized advice for refinement, enabling non-technical users to produce clearer and more complete test scenarios.

Our experiments show that small language models like CamemBERT and SmolLM2 successfully identify semantic similarities between natural language descriptions and their Gherkin translations while maintaining data privacy through local processing. Additionally, the interactive advice mechanism reduces the number of attempts required to meet quality standards by 26%, highlighting its effectiveness in user-driven improvements.

While GUIDE has shown promise in improving the quality of business-driven test descriptions, several avenues for improvement remain open. Notably, we observed issues related to the structuring of Gherkin translations during the evaluation process. Despite its structured format, Gherkin generated by the translation step sometimes suffers from syntactic inconsistencies or incorrect formatting, which can hinder the subsequent classification and assessment. To address this limitation, it could be possible to use syntax-aware models that validate Gherkin structure during generation, or to apply post-processing corrections to ensure compliance with Gherkin's strict syntax.

## Limitations

One of the core design choices of GUIDE is the use of small language models (CamemBERT and SmolLM2) to ensure local execution and respect for data privacy. While this choice enables on-premises deployment and reduces dependency on external cloud services, it also introduces a limitation in terms of generalization. Unlike larger pre-trained models (e.g., GPT-4, PaLM), smaller models require more task-specific fine-tuning to perform adequately. This additional training phase

can cause GUIDE to become more domain-specific, potentially limiting its effectiveness when exposed to new application contexts or unseen business-specific terminologies.

Moreover, GUIDE relies heavily on Gherkin as an intermediate representation to assess the quality of test descriptions. While Gherkin is well-structured and human-readable, it enforces a rigid format that may not capture more complex testing logic or non-linear interactions described by business experts.

# References

Saranya Alagarsamy, Chakkrit Tantithamthavorn, and Aldeida Aleti. 2024. A3test: Assertion-augmented automated test case generation. *Information and Software Technology*, 176:107565.

Loubna Ben Allal, Anton Lozhkov, Elie Bakouch, Gabriel Martín Blázquez, Guilherme Penedo, Lewis Tunstall, Andrés Marafioti, Hynek Kydlíček, Agustín Piqueres Lajarín, Vaibhav Srivastav, Joshua Lochner, Caleb Fahlgren, Xuan-Son Nguyen, Clémentine Fourrier, Ben Burtenshaw, Hugo Larcher, Haojun Zhao, Cyril Zakka, Mathieu Morlon, Colin Raffel, Leandro von Werra, and Thomas Wolf. 2025. Smollm2: When smol goes big – data-centric training of a small language model.

Wissam Antoun, Francis Kulumba, Rian Touchent, Éric de la Clergerie, Benoît Sagot, and Djamé Seddah. 2024. Camembert 2.0: A smarter french language model aged to perfection.

Jacob Austin, Augustus Odena, Maxwell Nye, Maarten Bosma, Henryk Michalewski, David Dohan, Ellen Jiang, Carrie Cai, Michael Terry, Quoc Le, and Charles Sutton. 2021. Program synthesis with large language models.

Zachary Eberhart and Collin McMillan. 2022. Generating clarifying questions for query refinement in source code search. In *2022 IEEE International Conference on Software Analysis, Evolution and Reengineering (SANER)*, pages 140–151. IEEE.

Yaojie Hu, Qiang Zhou, Qihong Chen, Xiaopeng Li, Linbo Liu, Dejiao Zhang, Amit Kachroo, Talha Oz, and Omer Tripp. 2025. QualityFlow: An agentic workflow for program synthesis controlled by LLM quality checks.

Albert Q. Jiang, Alexandre Sablayrolles, Arthur Mensch, Chris Bamford, Devendra Singh Chaplot, Diego de las Casas, Florian Bressand, Gianna Lengyel, Guillaume Lample, Lucile Saulnier, Lélio Renard Lavaud, Marie-Anne Lachaux, Pierre Stock, Teven Le Scao, Thibaut Lavril, Thomas Wang, Timothée Lacroix, and William El Sayed. 2023. Mistral 7b.

Carlos E Jimenez, John Yang, Alexander Wettig, Shunyu Yao, Kexin Pei, Ofir Press, and Karthik R Narasimhan. 2024. SWE-bench: Can language models resolve real-world github issues? In *The Twelfth International Conference on Learning Representations*.

Haau-Sing Li, Mohsen Mesgar, André FT Martins, and Iryna Gurevych. 2022. Python code generation by asking clarification questions. *arXiv preprint arXiv:2212.09885*.

Fangwen Mu, Lin Shi, Song Wang, Zhuohao Yu, Binquan Zhang, Chenxue Wang, Shichao Liu, and Qing Wang. 2023. Clarifygpt: Empowering llm-based code generation with intention clarification. *arXiv preprint arXiv:2310.10996*.

Olga Ormandjieva, Ishrar Hussain, and Leila Kosseim. 2007. Toward a text classification system for the quality assessment of software requirements written in natural language. In *Fourth international workshop on Software quality assurance: in conjunction with the 6th ESEC/FSE joint meeting*, pages 39–45.

Eugenio Parra, Christos Dimou, Juan Llorens, Valentín Moreno, and Anabel Fraga. 2015. A methodology for the classification of quality of requirements using machine learning techniques. *Information and Software Technology*, 67:180–195.

Reid Pryzant, Dan Iter, Jerry Li, Yin Tat Lee, Chenguang Zhu, and Michael Zeng. 2023. Automatic Prompt Optimization with "Gradient Descent" and Beam Search.

Jörg Tiedemann, Mikko Aulamo, Daria Bakshandaeva, Michele Boggia, Stig-Arne Grönroos, Tommi Nieminen, Alessandro Raganato Yves Scherrer, Raul Vazquez, and Sami Virpioja. 2023. Democratizing neural machine translation with OPUS-MT. *Language Resources and Evaluation*, pages 713–755.

Jörg Tiedemann and Santhosh Thottingal. 2020. OPUS-MT — Building open translation services for the World. In *Proceedings of the 22nd Annual Conferenc of the European Association for Machine Translation (EAMT)*, Lisbon, Portugal.

Michele Tufano, Dawn Drain, Alexey Svyatkovskiy, Shao Kun Deng, and Neel Sundaresan. 2021. Unit Test Case Generation with Transformers and Focal Context.

Jian Wang and Wenjie Li. 2021. Template-guided clarifying question generation for web search clarification. In *Proceedings of the 30th ACM international conference on information & knowledge management*, pages 3468–3472.

Sebastian Weigelt, Vanessa Steurer, and Walter F. Tichy. 2020. At your Command! An Empirical Study on How LaypersonsTeach Robots New Functions.

Zhuokui Xie, Yinghao Chen, Chen Zhi, Shuiguang Deng, and Jianwei Yin. 2023. ChatUniTest: a ChatGPT-based automated unit test generation tool.

Weimin Xiong, Yiwen Guo, and Hao Chen. 2023. The program testing ability of large language models for code. *arXiv preprint arXiv:2310.05727*.

J D Zamfrescu-Pereira, Richmond Wong, Bjoern Hartmann, and Qian Yang. 2023. Why Johnny Can't Prompt: How Non-AI Experts Try (and Fail) to Design LLM Prompts.

Fida Zubair, Maryam Al-Hitmi, and Cagatay Catal. 2024. The use of large language models for program repair. *Computer Standards & Interfaces*, page 103951.

## A  Similarity Scale

- **Completely Similar** (COMPSIM): it is the same test and the same procedure expressed a little differently.

- **Very Similar** (VERSIM): it is the same test, but the operating procedures are different (one may be more detailed than the other, but that does not mean that the test described is different).it is the same test, but the operating procedures are slightly different.

- **Somewhat Similar** (SOMSIM): it is probably the same test (but I am not sure) and/or there are many differences in the operating methods used.

- **Very Different** (VERDIFF): it may be the same test and/or the operating procedures expressed have too many differences (but elements in common).

- **Completely Different** (COMPDIFF): it is not the same test and/or the operating procedures are completely different (no common elements).

## B  Esco Explorer



Figure 3: Screenshot of the software Esco Explorer

## C  Prompt for the Beam Search

All the prompts presented in this appendix were originally written in French.

Listing 1: The prompt used as the basis for the Beam Search.

```
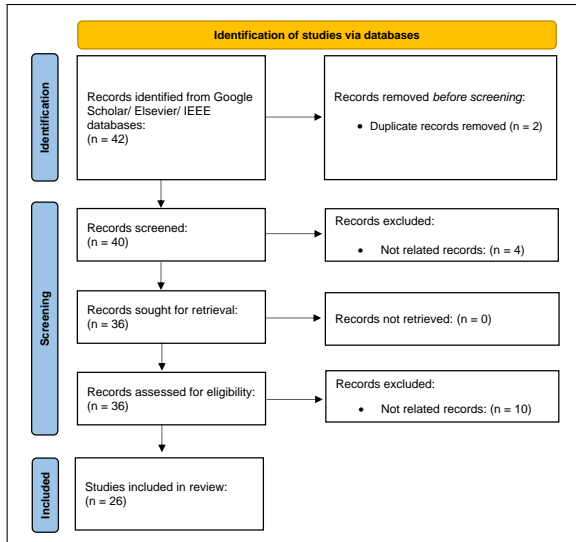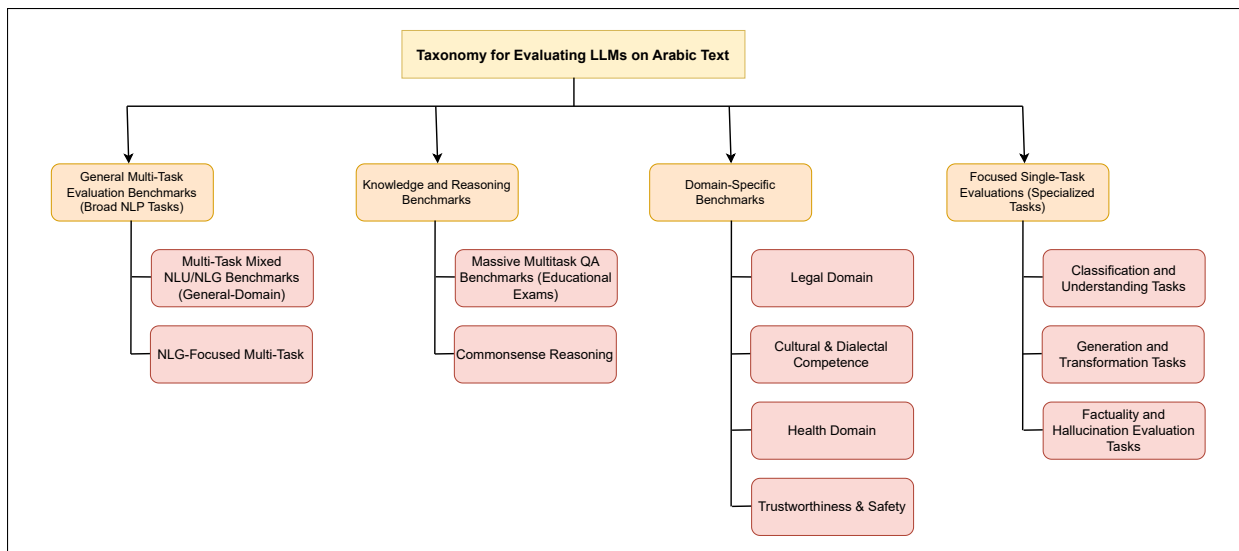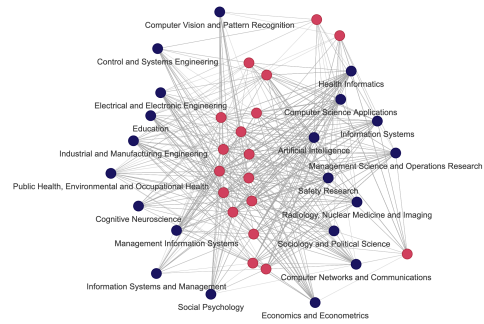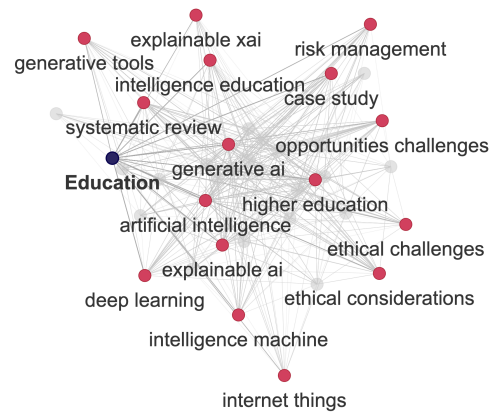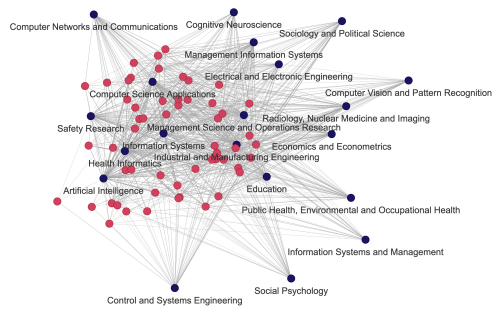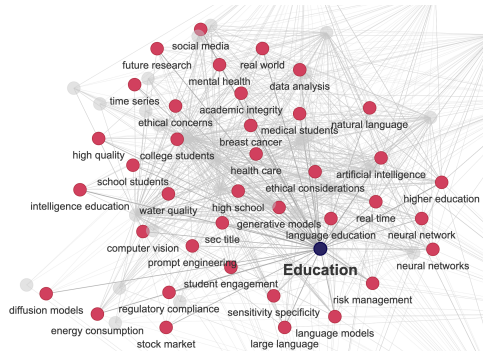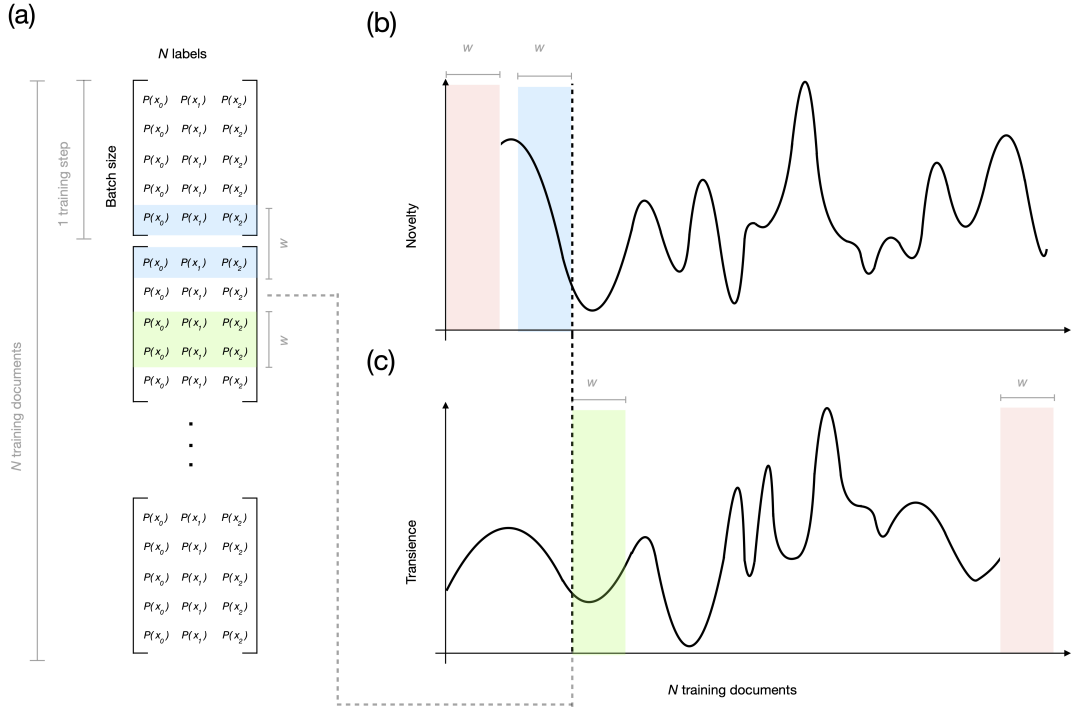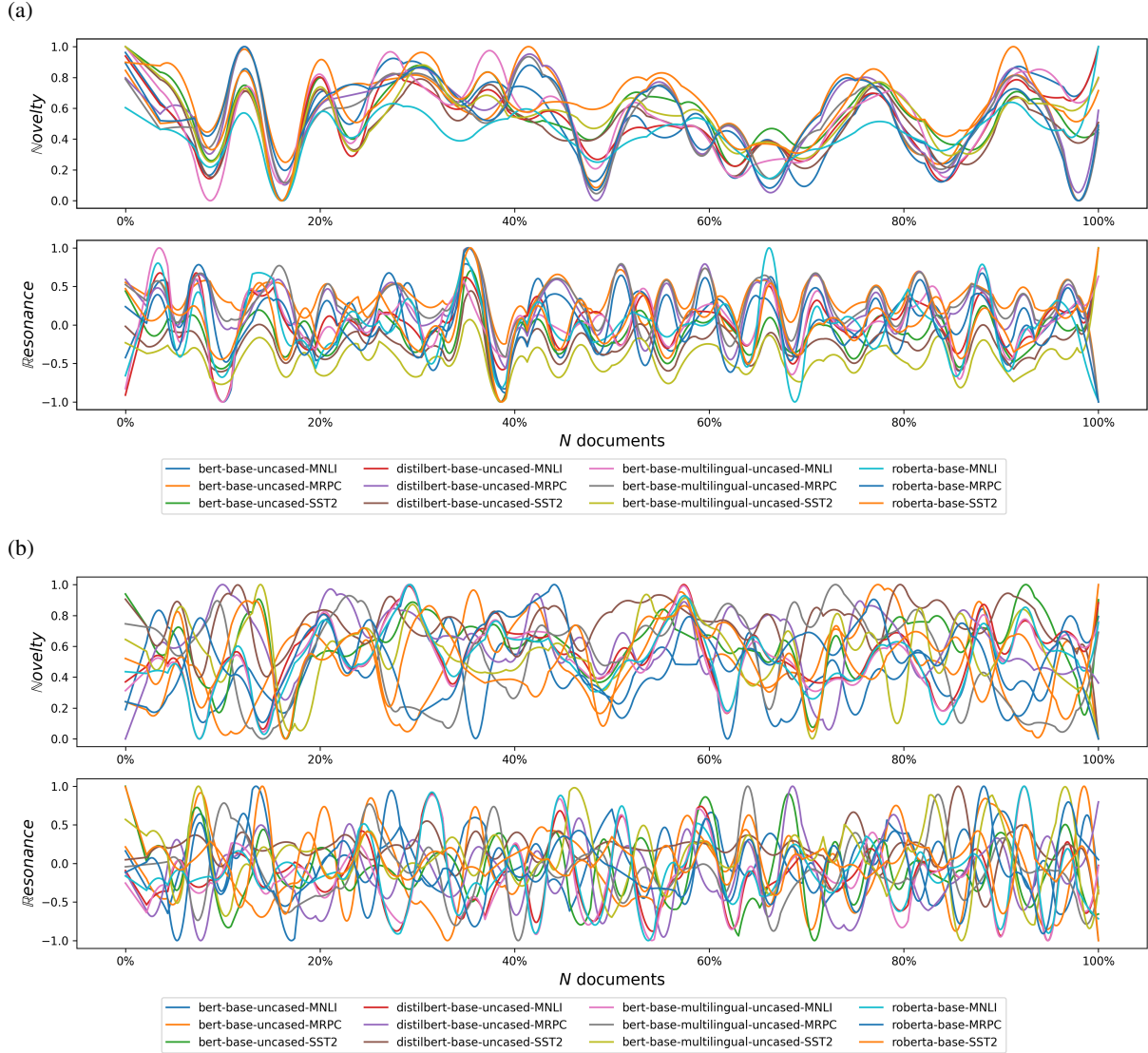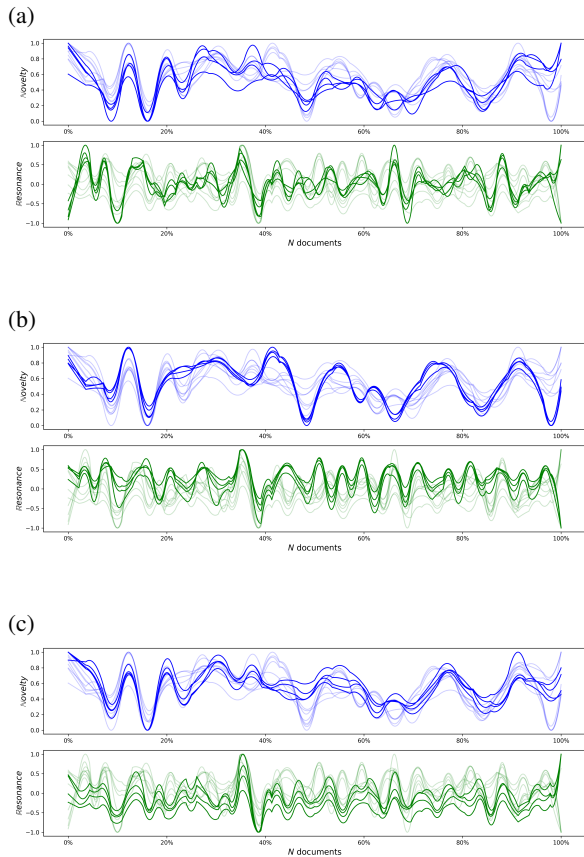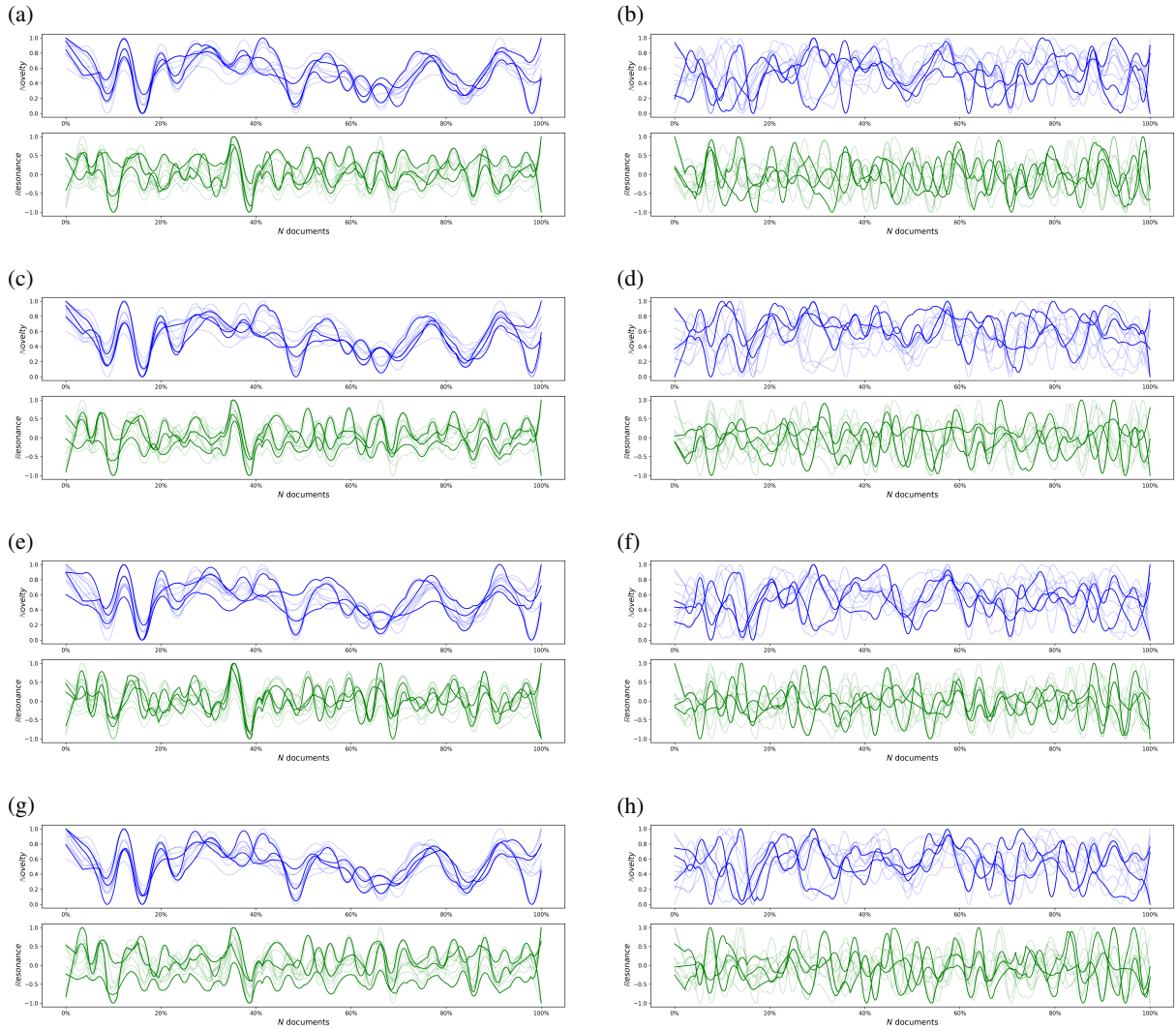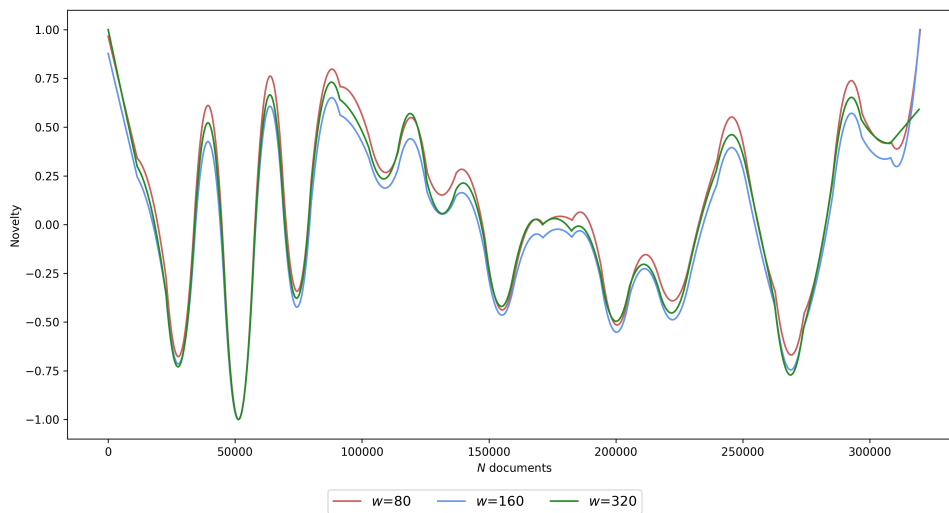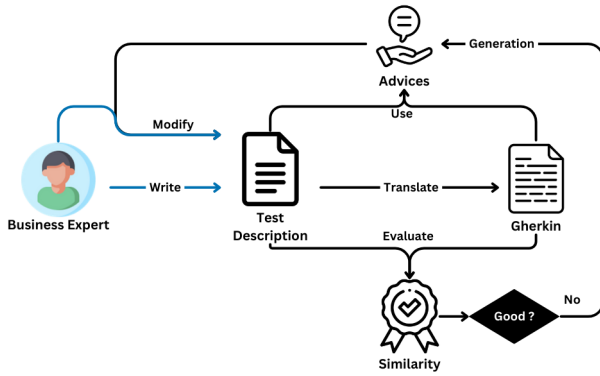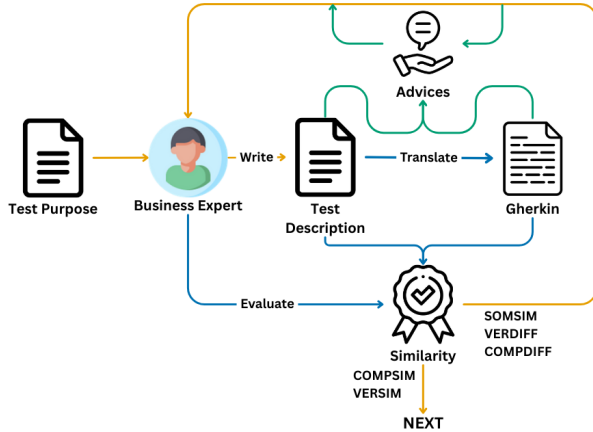The user wrote a test description , and an automatic
analysis identified the elements that were missing
or needed to be improved.
Based on the following information:
- Test description provided:
"%s"

- Gherkin:
"%s"

- Software context:
"%s"

Generate clear, actionable suggestions to help the
user improve their description. Your recommendations
 must be precise and adapted to the elements
detected as insufficient or missing. Respond only to
 suggestions and nothing else.
```

Listing 2: The prompt used to make the mutation during the Beam Search.

```
Generates three variations of the following
instruction, while retaining its semantic meaning.
Each variation must respect the following
constraints:

The message must remain clear and understandable.
The three %s markers must be retained, the first for
 the description, the second for the Gherkin and the
 third for the software.
The variations must reformulate the instruction
without changing its content.
Separate each variation with "---".

Reference instruction:
"%s"
```

Listing 3: The prompt considered to be the best during the Beam Search.

```
The user has drawn up a test description , and an
automatic analysis has highlighted elements to be
enhanced or corrected.

Based on the following information:
- Test description provided:
"%s"

- Gherkin:
"%s

- Software context:
"%s

Make concrete, achievable suggestions to help the
user improve his description. Your recommendations
must be explicit and adapted to the elements deemed
to be missing or insufficient. Answer only the
suggestions and nothing else.
```

# D   Data

Table 6: Examples illustrating the disparity of functional test descriptions produced by different users. These descriptions, taken from the production campaign, all aim to express the same test objective: 'File - Open a CSV file'. The similarity label was assigned by the user who wrote the description

| Description | Gherkin | Similarity Label |
| --- | --- | --- |
| Cliquer sur "OPEN". Vérifier que la fenêtre de choix des fichiers s'ouvre. Sélectionner un fichier avec l'extension ".csv" et vérifier que la recherche apparaît dans le logiciel. | **Given** le logiciel est ouvert<br>**When** je clique sur "OPEN"<br>**Then** la fenêtre de choix des fichiers s'ouvre<br>**And** je sélectionne un fichier avec l'extension ".csv"<br>**And** la recherche apparaît dans le logiciel | VERSIM |
| Dans ESCO Explorer, appuyer sur le bouton "OPEN" tout en haut à gauche. Une fois fait, chercher dans la liste un fichier avec comme type de fichier "Fichier CSV Microsoft Excel". Double cliquer sur le nom du fichier pour l'ouvrir. Si le fichier s'ouvre dans ESCO, on peut dire que le test est validé. | **Given** l'utilisateur est sur l'application ESCO.<br>**When** l'utilisateur sélectionne un fichier CSV via le bouton "OPEN".<br>**Then** le fichier CSV choisi devrait s'ouvrir dans ESCO. | COMPSIM |
| Cliquer sur le bouton 'OPEN' de la barre de menu d'EscoExplorer. Cliquer ensuite sur le fichier nommé 'cobolview.csv' dans la fenetre qui s'est ouverte puis cliquer sur le bouton 'ouvrir'. Le test est réussi si et seulement si le mot 'COBOL' apparaît dans la fenetre de recherche d'EscoExplorer. | **Given** EscoExplorer a été lancé<br>**And** le langage sélectionné est 'English'<br>**When** cliquer sur le bouton 'OPEN'<br>**And** ouvrir le fichier 'cobolview.csv'<br>**Then** 'COBOL' apparait dans la fenetre de recherche | SOMSIM |

# E  Campaign Analysis

Table 7: Average number of productions per test category, dividing the test categories into two halves.

| Test Category | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 |
|---|---|---|---|---|---|---|---|---|---|
| **First Half** | 2.97 | 1.5 | 2.6 | 1.45 | 1.8 | 1.52 | 1.62 | 1.63 | 1.5 |
| **Second Half** | 2.02 | 1.6 | 1.72 | 1.9 | 1.47 | 1.63 | 1.5 | 1.62 | 1.43 |

Table 8: User Behavior During the Production Campaign (Number of attempts to meet quality criterion)

| User ID | Avg. Attempts | 1 Attempt | 2 Attempts | 3 Attempts | 3+ Attempts | Abandonment |
|---|---|---|---|---|---|---|
| **0** | 1.71 | 42 | 12 | 7 | 4 | 3 |
| **1** | 1.15 | 59 | 8 | 1 | – | – |
| **2** | 3.07 | 18 | 18 | 12 | 20 | – |
| **3** | 1.91 | 32 | 19 | 12 | 5 | – |
| **4** | 2.31 | 21 | 11 | 10 | 1 | 25 |
| **5** | 1.79 | 35 | 18 | 8 | 5 | 2 |
| **6** | 1.29 | 53 | 12 | 2 | 1 | – |
| **7** | 1.38 | 51 | 13 | 2 | 2 | – |
| **8** | 1.44 | 45 | 16 | 6 | – | 1 |
| **9** | 1.41 | 46 | 12 | 5 | – | 5 |
| **Total** | 1.75 | 402 | 139 | 65 | 38 | 36 |

Table 9: User Behavior During the Advice Campaign (Number of attempts to meet quality criterion)

| User ID | Avg. Attempts | 1 Attempt | 2 Attempts | 3 Attempts | 3+ Attempts | Abandonment |
|---|---|---|---|---|---|---|
| **10** | 1.21 | 60 | 6 | – | 2 | – |
| **11** | 1.06 | 64 | 4 | – | – | 1 |
| **12** | 1.19 | 59 | 6 | 2 | 1 | 1 |
| **13** | 1.16 | 58 | 9 | 1 | – | – |
| **14** | 1.15 | 59 | 8 | 1 | – | – |
| **15** | 1.15 | 49 | 11 | 4 | 4 | 3 |
| **16** | 1.18 | 60 | 4 | 4 | – | – |
| **17** | 1.85 | 36 | 8 | 22 | 2 | 18 |
| **Total** | 1.29 | 445 | 46 | 34 | 9 | 23 |

# Tokenization and Morphology in Multilingual Language Models: A Comparative Analysis of mT5 and ByT5

**Thao Anh Dang**
Utrecht University
Radboud University
t.t.a.dang@uu.nl

**Limor Raviv**
Max Planck Institute
for Psycholinguistics
limor.raviv@mpi.nl

**Lukas Galke**
University of Southern Denmark
galke@imada.sdu.dk

## Abstract

Morphology is a crucial factor for multilingual language modeling as it poses direct challenges for tokenization. Here, we seek to understand how tokenization influences the morphological knowledge encoded in multilingual language models. Specifically, we capture the impact of tokenization by contrasting a minimal pair of multilingual language models: mT5 and ByT5. The two models share the same architecture, training objective, and training data and only differ in their tokenization strategies: subword tokenization vs. character-level tokenization. Probing the morphological knowledge encoded in these models on four tasks and 17 languages, our analyses show that the models learn the morphological systems of some languages better than others and that morphological information is encoded in the middle and late layers. Finally, we show that languages with more irregularities benefit more from having a higher share of the pre-training data.

## 1 Introduction

Tokenization, the process of segmenting a text into individual units, plays a special role in language modeling as it is disconnected from the otherwise end-to-end training procedure (Xue et al., 2021, 2022; Sennrich et al., 2016). Languages differ in their morphological structure (Goldman and Tsarfaty, 2022; Dryer and Haspelmath, 2013; Evans and Levinson, 2009; Ackerman and Malouf, 2013) and it has been shown that morphologically more complex languages are harder to acquire by humans (DeKeyser, 2005; Raviv et al., 2021; Kempe and Brooks, 2008) and deep neural network models (Galke et al., 2024; Park et al., 2021; Mielke et al., 2019; Cotterell et al., 2018). Here, we seek to understand to what extent different tokenization strategies influence the ability of multilingual language models to capture morphological knowledge in different languages (See Figure 1).

Ideally, a language model would be equally proficient in a variety of languages (Lample and Conneau, 2019; Conneau et al., 2020; Ruder et al., 2019). Understanding the influence of tokenization is crucial in the context of multilingual language modeling (Xue et al., 2021, 2022; Warstadt et al., 2020), as it is challenging to find a set of tokens that is equally good for modeling all the languages in the world. Beyond the proportions of languages in the training data, it is important to take into account the morphological structure of the different languages (Anh et al., 2024; Galke et al., 2024; Cotterell et al., 2018). Importantly this needs to be already considered when selecting the data for learning the tokenizer – even before language model pre-training – as this influences what subword structures end up as the tokens to be processed by the language model. The issue of tokenization and, in related matter, how to mix different languages, are particularly relevant in the current era of large language models (Touvron et al., 2023; Brown et al., 2020; Bubeck et al., 2023; Wei et al., 2022), when aiming for similar performance on a diverse range of languages (Le Scao et al., 2023).

With models such as ByT5 (Xue et al., 2022), a character-level language model based on the T5 architecture (Raffel et al., 2020), it has been shown that tokenizer-free language models yield commensurate downstream performance with their tokenizer-based counterparts (Xue et al., 2022; Edman et al., 2024), such as mT5 (Xue et al., 2021), another T5-based model that trained on the exact same data as ByT5. Specifically, in machine translation, Edman et al. (2024) have found that character-level ByT5 yields similar performance as mT5 when allowing more training to recover word-level structures. Yet, the interplay of morphology and tokenization is so far poorly understood.

Here, we seek to dissect the root of these findings through analyzing the effect of the tokenization strategy (character-level vs. subword-level) on the

Figure 1: Overview of our experimental procedure

morphological knowledge encoded in contextualized representations of multilingual language models. Specifically, we contrast ByT5 as a character-level multilingual language model and mT5 as a subword-level multilingual language model, with both models sharing the same architecture and being trained on the same data. We use well-established structural probing techniques to capture the amount of linguistic information encoded in the contextualized word representations of the language models (Rogers et al., 2021; Manning et al., 2020; Belinkov et al., 2020; Tenney et al., 2019). Focusing on morphology, we probe the contextualized representations of ByT5 and mT5 for morphological knowledge on 17 languages from a large-scale multilingual dataset (Acs et al., 2023). In addition, we use the non-contextualized fastText model (Bojanowski et al., 2017) as a control to understand to what extent the context is important for multilingual language models reflecting morphosyntactic structures in their representations. We further explore to what extent the captured morphological knowledge depends on the share that the languages have in the pre-training data, various linguistic factors, such as the type of task (number, tense, case, gender), and the language's degree of morphological complexity, as quantified by the degree of irregularity (Wu et al., 2019) and type-to-token ratio (TTR) (Bentz et al., 2015).

By systematically contrasting tokenizer-free ByT5 and subword-tokenized mT5, two pre-trained multilingual language models based on the T5 architecture, and trained on the same data, we find:

- Multilingual language models learn the morphological systems of some languages better than others.

- Morphological knowledge representation improves over transformer layers.

- Both subword- and byte-level models display approximate the same level morphological knowledge encoded in the activations after a small number of transformer layers.

- A language's degree of irregularity plays a substantial role for capturing morphological knowledge, suggesting that more irregular languages would benefit from a higher proportion of training data.

## 2 Background and Related Work

**Morphology** Morphology concerns how meaningful word units can be combined to express a range of grammatical information (Bloomfield, 1933). Languages differ greatly in their morphological systems and degrees of morphological complexity (Dryer and Haspelmath, 2013; Lupyan and Dale, 2010; Bloomfield, 1933). It has been shown that morphologically more complex languages are harder to acquire by humans (DeKeyser, 2005; Raviv et al., 2021; Kempe and Brooks, 2008) and deep neural networks (Galke et al., 2024; Cotterell et al., 2018). Some studies specifically investigate the link between morphological complexity and the challenges in language modeling (Cotterell et al., 2018; Mielke et al., 2019; Park et al., 2021; Galke

et al., 2024; Anh et al., 2024). However, the findings have been mixed. As most of these studies focus more on the overall learnability (Park et al., 2021; Gerz et al., 2018), there may be confounds other than morphological complexity, which we isolate here with a fine-grained analysis.

**Linguistic Probing**  Linguistic probing can be categorized into behavioral probing and structural probing (Madsen et al., 2022). Behavioral probing aims to understand how a language model behaves in a specific or new setting (Hupkes et al., 2023). Structural probing instead seeks to localize where linguistic abilities are encoded (Rogers et al., 2021). It was suggested that morphological knowledge is mainly encoded at the lower layers (Peters et al., 2018; Belinkov et al., 2020). However, most studies focus on analyzing monolingual language models, such as BERT (Devlin et al., 2019), or neural machine translation systems trained on pairs of languages (Acs et al., 2023; Belinkov et al., 2020; Bisazza and Tump, 2018; Edmiston, 2020), while studies on multilingual language models are rare.

**Tokenization**  Tokenization is a crucial step in language modeling, especially for multilingual language models. The dominant tokenization methods for language models are based on the Byte-Pair Encoding (BPE) algorithm (Sennrich et al., 2016), and rely on the data mix to ensure coverage of different languages (Le Scao et al., 2023). Starting with single characters, BPE iteratively merges tokens based on co-occurrence statistics, allowing for both subword and multi-word tokens. Comparing BPE with other subword tokenizers, Ali et al. (2024) found that the preference for tokenization methods differs across languages. While BPE works better for Germanic languages, such as German and English, Unigram (Kudo, 2018) is more well-suited for Romance languages, such as Spanish. The importance of tokenization in multilingual language modeling is evident (Hofmann et al., 2021; Toporkov and Agerri, 2024), as many studies propose new tokenization algorithms to make language models capture the morphological structure of the different languages (Goldman and Tsarfaty, 2022).

**Character-level Language Models**  A line of research investigates character-level language models to circumvent the issues of multi-lingual tokenization (Fleshman and Van Durme, 2023; Clark et al., 2022; Xue et al., 2022; Gao et al., 2020; Chung et al., 2016; Lee et al., 2017; Kim et al., 2016). In

machine translation, Lee et al. (2017) found that character-level tokenizers perform as well as or better than models based on sub-word tokenization. They highlighted that character-level tokenization offers better translation quality in multilingual and low-resource settings because of the shared vocabulary. Edman et al. (2024) argued that character-level models are better at learning information that operates at a low level of granularity, such as morphology. Comparing translation capability between multilingual language models with standard tokenization (mT5) and character-level model (ByT5), they found that ByT5 outperforms mT5 in several aspects. First, ByT5 produces higher-quality translations than mT5, even in the case of low-resource languages. ByT5 is also better in handling rare and similar words. In addition, low-resource languages may benefit from shared vocabulary, namely the set of characters (Gao et al., 2020). However, it has also been suggested that the effect varies across languages (Ali et al., 2024), which in-turn motivates the present study.

**Summary**  Overall, earlier work on probing the morphological knowledge of language models is mainly conducted on machine translation and recurrent neural networks (Belinkov et al., 2020; Vylomova et al., 2017). More recently, the focus has been extended to Transformer-based language models (Acs et al., 2023; Edmiston, 2020; Wu and Dredze, 2020), yet the majority of studies investigated BERT and its variants (Rogers et al., 2021). Previous studies provide mixed findings about the morphological knowledge of language models and only very few multilingual studies. Here, we complement the literature by analyzing how captured morphological knowledge is influenced by the tokenization strategy, the languages' proportions in the model's training data, the language's morphological complexity, and the task type.

## 3 Models and Data

We compare two pre-trained language models: mT5 and ByT5. The two models share the same architecture and are trained on the same data with the same training objective (masked span prediction). The key difference between mT5 and ByT5 is their tokenization strategy: mT5 uses a standard subword tokenization strategy, whereas ByT5 operates on character level. Thus, we can investigate the effect of tokenization.

**mT5** The mT5 model (Xue et al., 2021) is a multilingual version of T5, an encoder-decoder language model (Raffel et al., 2020). It is trained on the mC4 corpus (Xue et al., 2021), which consists of text in 101 languages compiled from the Common Crawl web scrape. Like T5, mT5 employs a masked span prediction training objective and the SentencePiece tokenizer (Kudo and Richardson, 2018), a variant of the BPE tokenizer by Sennrich et al. (2016). The vocabulary comprises approximately 250,000 subwords, covering 104 languages (Xue et al., 2021) while sharing (sub-)words between languages.

**ByT5** ByT5 (Xue et al., 2022) is a tokenizer-free variant of mT5, inheriting most of the properties of mT5, using the same T5 model architecture, and the same masked span prediction training objective. The only crucial difference between them is the tokenization method: While mT5 uses a standard tokenizer, ByT5 operates on single-character tokens, or more precise: UTF-8 encoded bytes. Another difference is that ByT5's encoder stack consists of three times more layers than the decoder to process a larger number of bytes. For the masked span prediction objective, ByT5 uses a span of 20 bytes. The similarity in architecture and sizes of mT5 and ByT5 enables our comparison of how the tokenization strategy impacts the morphological knowledge encoded in the learned representations.

**Dataset** We use the multilingual morphological probing dataset by Acs et al. (2023). The dataset consists of 247 probing tasks, available in 42 languages and 10 language families. It is built upon the Universal Dependencies tree bank and covers both frequent and infrequent words. In each language, each task includes a training set of 2000 examples, a development test of 200 examples, and a test set of 200 examples. We selected 16 out of all 42 languages which had at least two tasks available. However, we also included Arabic despite having only one task available to better cover the Semitic language family. In total, our considered dataset consists of 17 languages and 43 morphological probing tasks, covering number, case, gender, and tense – focusing on nominal and verbal inflection.

## 4 Methodology

**Feature Extraction** For training, we extracted the contextualized embeddings of the words in the training set for each task in each language and each model. Both mT5 and ByT5 are available

in different sizes. We chose to test the `mT5-base` model. We froze the weights and extracted the hidden states for the entire sentence before extracting the word embedding corresponding to the target. Since we also aim to look at how much morphological knowledge is learned at each layer, we extracted the word embedding at each hidden layer of the network, including the input embedding layer. Each word representation is associated with a label, which corresponds to the respective morphological feature from the task. We trained separate probes for each language and each task in each layer of the model.

**Probing Classifiers** Previous studies often use two architectures for probing classifiers, namely linear classifiers (Hupkes et al., 2018; Belinkov et al., 2020) and multilayer perceptrons (MLPs) (Lin et al., 2019; Conneau et al., 2018; Adi et al., 2017; Ettinger et al., 2018; Zhang and Bowman, 2018). Both types of probes have received convincing arguments. Linear probes capture information that is linearly separable in the representations (Liu et al., 2019; Belinkov et al., 2020), whereas MLP probes can additionally capture nonlinear patterns in the representations (Hewitt and Liang, 2019). Studies show that linear classifiers and MLPs produce similar accuracy (Conneau et al., 2018; Belinkov et al., 2017a; Qian et al., 2016). We have considered both types of probes for this study (see Appendix C). For the main results, we employ MLP probes.

**Subword Pooling** In both mT5 and ByT5, words are segmented into either subword units or characters. As such, when passing through the hidden layers, each subword or character has its own embedding. There are several methods to then approximate the embedding for an entire word: The first method is to take the weighted **average** of the embeddings of all components. The second way is to consider the embedding of the **last** subword or character as the representation for the entire word. Both methods have limitations. Averaging the token embeddings may cancel out some information, while the last embedding may not contain all the information about the entire word. Belinkov et al. (2020) compared both methods and found that using the embeddings of the last token produced higher accuracy scores. We have considered both options (see Appendix D) and can confirm that last-token pooling leads to higher probing accuracy. For the main results, we employ last-token pooling.

**Evaluation and Control** We aim to probe the morphological knowledge of a range of typologically diverse languages. While the dataset of Acs et al. (2023) supports 42 languages, in some languages, there is a large gap between the number of tasks in each language. While Russian has 12 tasks, Polish and Armenian have only one task. To ensure a fair comparison of the learned morphological representations, we selected languages with at least two tasks along with Arabic to cover the Semitic language family. We focused exclusively on the morphological properties of words, excluding agreement tasks. In total, we ran 43 morphological probing tasks for 17 languages. Appendix E provides the details of the morphological properties of the 17 considered languages, their morphological complexity scores, and their proportion in the ByT5/mT5 training data. Appendix F provides a detailed description of the types of probing tasks, covering number, case, gender, and tense.

As non-contextualized control, we employ fast-Text word embeddings (Bojanowski et al., 2017) which are available in 157 languages. We probed fastText embeddings using the exact same procedure and evaluation metric. We obtained the static word embeddings without any pooling. Contrasting the contextualized representations of mT5 and ByT5 with fastText allows us to quantify to what extend the contextualized models make use the sentence context to tackle the morphological tasks.

## 5 Results

**Overall Probing Accuracy** We first looked at the overall probing performance of mT5, ByT5, and fastText as well as the differences between the two Transformer-based models compared to the fastText baseline. For all analyses except for the layer-wise analysis, we used the probing accuracy of the last hidden layer. To obtain the overall performance of mT5 and ByT5, we averaged over all languages and tasks, resulting in a single accuracy score for each model (see Table 1). Full results for each task and language are provided in Appendix G.

On the surface, it appears that mT5 and ByT5 have comparable performance and both models outperformed fastText. ByT5 slightly surpassed mT5, yet this difference is very small. This finding is different from that of Belinkov et al. (2017a), who found that character-level tokenizers are better than subword tokenizers in representing morphology.

| Model | Mean Probing Accuracy |
|---|---|
| mT5-base | 82.57 |
| ByT5-base | **82.86** |
| fastText (baseline) | 77.52 |

Table 1: Probing accuracy of mT5, ByT5 and fastText, averaged over languages and tasks

Table 2 shows the difference between accuracy scores of mT5 and ByT5, grouped by language. It can be seen that accuracy scores are not equal across both languages and tasks. Comparing mT5 and ByT5, it seems that they perform on par with each other in most languages. However, mT5 scores higher in Turkish and ByT5 achieves much higher on Hindi tasks. The difference between contextualized language models and non-contextualized fastText also tells to what extent contextual information affects morphological abilities. From the results, we observed that for most languages, mT5 and ByT5 achieved considerably higher probing accuracy than the baseline. The largest difference is observed in the case of Basque (> 50%). This implies that contextual word embeddings capture morphological knowledge better than static embeddings for these languages. In contrast, the results are lower than the non-contextual baseline in French and Russian. Contextual information seems to make accessing morphological information more difficult in these two languages.

Considering the differences between both language models and the fastText baseline, as shown in Table 2, it can be seen that they generally outperformed the baseline yet perform substantially worse than baseline in French and Russian. Averaging over tasks, ByT5 yields the highest probing accuracy in 7 languages, while mT5 yields the highest probing accuracy on 6 tasks out of 17 languages.

Table 3 shows the probing accuracy scores for each language, averaged over tasks. It appears that there are some differences in accuracy across languages, with the hardest language being Russian for mT5 and Arabic for ByT5. The accuracy scores of some languages are higher than the others. Unsurprisingly, the models perform best on the English language, followed by Hebrew, Portuguese, and Romanian. The models learn the morphological systems of German, French, Estonian, and Latvian moderately well. Arabic and Russian achieved lowest accuracy scores. However, Arabic results should be interpreted with caution as there

| Language | Family | Model | | |
|---|---|---|---|---|
| | | mT5 | ByT5 | fastText |
| English | Germanic | 98.00 | **98.50** | 97.75 |
| Dutch | Germanic | **93.50** | 92.75 | 82.25 |
| German | Germanic | 68.46 | **72.40** | 68.87 |
| French | Romance | 71.93 | 77.53 | **92.95** |
| Spanish | Romance | 93.83 | **94.00** | 71.16 |
| Portuguese | Romance | 95.00 | **96.25** | 88.16 |
| Romanian | Romance | 94.75 | **95.25** | 92.25 |
| Hebrew | Semitic | 97.50 | **98.00** | 92.49 |
| Arabic | Semitic | **66.17** | 49.25 | 37.81 |
| Russian | Slavic | 61.26 | 57.43 | **79.20** |
| Czech | Slavic | **88.79** | 87.79 | 78.81 |
| Hindi | Indic | 67.83 | **87.50** | 58.02 |
| Urdu | Indic | **88.50** | 84.50 | 74.33 |
| Turkish | Turkic | **94.78** | 83.69 | 78.28 |
| Latvian | Baltic | **77.14** | 72.01 | 73.76 |
| Estonian | Uralic | **91.08** | 90.03 | 82.94 |
| Basque | Basque | 91.69 | **91.91** | 52.60 |

Table 2: Probing accuracy of mT5, ByT5, and fastText by languages, language families, averaged over tasks. The best score per language is marked in bold font.

| Task | mT5 | ByT5 | fastText |
|---|---|---|---|
| Number | **93.75** | 93.56 | 88.15 |
| Tense | 80.30 | 72.44 | **80.89** |
| Gender | 76.85 | **85.16** | 77.61 |
| Case | **71.53** | 69.16 | 55.49 |

Table 3: Probing accuracy of mT5, ByT5, and fastText, averaged over languages and tasks

is only one task available (i.e., case).

**Layer-wise Analysis** Figure 2 illustrates the difference in probing accuracy between languages and between tasks for mT5 and ByT5. It can be seen that there are some degrees of variation between layers. In languages that show high overall performance, namely English, Dutch, Portuguese, Spanish, Basque, and Hebrew, probing accuracy shows very little improvement over layers. In other languages, accuracy increases, reaches its peak at the middle and slightly decreases at late layers in other languages. This finding is partly consistent with Acs et al. (2023), Edmiston (2020), and Hewitt et al. (2021), who also reported best performance in the middle to late layers. However, we further show that this is not true for all languages. There are cases where morphological knowledge is successfully learned in the early layer and carried on throughout the network. Our results contrast with

Belinkov et al. (2020) and Peters et al. (2018), who found that morphological information is best encoded in the first layer of the model, and then has the tendency to decrease over time.

Comparing mT5 and ByT5, there are a few noticeable differences. In the plots for ByT5, accuracy scores of each language and each task improves considerably after the embedding layer. This trend is less visible in mT5, although performance does improve over layers. Morphology is better learned in the embedding layer of mT5 than that of ByT5. This may imply that character-level language models need more layers to capture morphological patterns of languages.

**Effect of Task Type** To investigate whether morphological features are learned differently by mT5 and ByT5, we averaged the scores for each task across all languages, resulting in a single score for each task (see Table 3). The results strongly suggest that each morphological feature is encoded differently. Interestingly, mT5 and ByT5 show different patterns. Both models perform equally well at number and worst at case. However, tense is learned better than gender by mT5 while the opposite is true for ByT5. Comparing both language models with the baseline, they surpass the baseline in all tasks except for tense, where ByT5 performs considerably worse than fastText.

Case seems to be the hardest task for both models. Besides the case task often having more possible classes than other features, case is also more context-dependent than other features. Case marking is used to indicate the syntactic function of the word in the sentence. As such, one word may have different cases in different contexts and thus is inflected distinctively. Gender is also relatively difficult, especially for mT5. However, looking at individual languages (see Table 2), the mean score of mT5 is affected by Hindi and Latvian, whose scores are exceptionally lower than the baseline (less than 25%). Except for those two languages, mT5 and ByT5 perform equally well.

**Morphological Complexity and Training Data** We explore the effects of two types of morphological complexity, namely TTR (Bentz et al., 2015) and the degree of irregularity (Wu et al., 2019) on probing accuracy. Complexity values per language can be found in Appendix E. We hypothesize that probing accuracy is influenced by the proportion of the respective language in the training data, which then also modulates the effect of a language's mor-

Figure 2: Probing accuracy of ByT5 (left) and mT5 (right) across layers grouped by languages and tasks. Each line represents a language (top) or a task (bottom). Each data point is the accuracy scores at each layer of each language.

phological complexity. To understand the effect of these factors and their interaction, we fitted a generalized mixed effect logistic regression model predicting accuracy from the two morphological complexity measures and the proportion of training data. There are 200 data points for each task in each language. The analysis was conducted with the R-package lme4. All variables were scaled and centered. Random intercepts were added for language and task.

Our results show significant positive effects of training data size and irregularity on probing accuracy, as well as their interaction. In more detail, there is a strong effect of the amount of training data and a language's degree of irregularity on probing accuracy (training data: $\beta = 2.11$, $SE = .55$, $p < .001$, irregularity: $\beta = 2.83$, $SE = .42$, $p < .001$). The effects of training data and irregularity were highly correlated (0.831). In addition, there is an interaction effect between the language's irregularity and its training data in the model ($\beta = 2.03$, $SE = .31$, $p < .001$). The effect of training data size on probing accuracy is stronger when there is more irregularity in the language. This means that high irregularity in the morphological systems amplifies the impact of training data on the morphological abilities of language models. Detailed results of the statistical models can be found in Appendix H.

## 6 Discussion

**Languages' morphology is learned differently**
Our probing results in mT5 and ByT5 show that the morphological knowledge of some languages is better represented than the others. Some languages (e.g., English, Dutch) achieved nearly perfect accuracy in probing tasks (higher than 90%). However, both mT5 and ByT5 performed worse at German, French, Russian, and Arabic tasks. These results to some extent contradict Edmiston (2020), who found comparable performance in all languages. Acs et al. (2023) also did not observe performance differences across languages, but only between part-of-speech and morphological features.

**Differences across tasks** We observed that some tasks are more difficult for the language models: Number is the easiest task whereas case is the hardest one. This difference can be partly attributed to the higher number of possible categories but also to the context-dependent nature of case. This is supported by the non-contextualized baseline results for case, which are substantially lower than both mT5 and ByT5. These findings are in agreement with earlier findings by Bisazza and Tump (2018) and Edmiston (2020). Why tense and also case features are particularly challenging to find in the representations of character-based language models is an interesting question for future research.

**Character-level models are on par with subword level models in representing morphology** We found that both models yield highly similar average performance on the probing tasks. Remarkably, despite being tokenized at the byte level, ByT5 is able to reconstruct morphological knowledge already in the first transformer block, arriving at a similar level as mT5. Belinkov et al. (2020) and Vylomova et al. (2017) tested machine translation systems and found character-level tokenizers to surpass BPE in learning morphology. However, through investigating large-scale language models, we found no advantage of character-level tokenizer in encoding morphology. Instead, we found that the preferable model differs between individual languages, e.g., subword tokenization works better for Turkish, whereas character-level tokenization benefits Hindi morphology. fastText performed nearly as well as mT5 and ByT5. We attribute this to the word-level nature of morphological features. Moreover, we used separated fastText models for each language, instead of generic multilingual ones.

**Morphology is best represented in middle to late layers** Our findings show that morphological knowledge generally improves over layers in both language models. There are languages which have high performance across layers. Yet, ByT5 shows greater improvement after the embedding layer than mT5. We also observed that morphological information is best encoded in the middle to late layers of the models in some languages. This finding supports Acs et al. (2023), Hewitt et al. (2021), and Edmiston (2020). Our findings differ from Belinkov et al. (2017b); Tenney et al. (2019); Peters et al. (2018), who found that morphology is a low-level feature and is encoded along with word identity in the first layer of the network.

**Morphological irregularity amplifies the effect of training data** Considering the relationship between morphological knowledge encoded in language models and the languages' morphological complexity, our analysis reveals effects of morphological irregularity and training data sizes on the performance of probing classifiers, in a way that the effect of irregularity is mediated by training data. When a language is highly irregular, a larger share of the training data is beneficial to fully capture its morphological system. Previous studies on the effect of training data sizes show that its effect is not present at the representation level yet at the downstream level (Warstadt et al., 2020; Zhang and Bowman, 2018). We have shown here that the importance of the relative training data size in multilingual language modeling can be already found when probing for morphological knowledge. Considering the interplay with morphological complexity, Mielke et al. (2019) and Gerz et al. (2018) correlated modeling difficulty with morphological counting complexity (Sagot, 2013), vocabulary sizes of languages, and dependency length – and found vocabulary size to be the most important factor. Our study complements those findings by establishing that the degree of irregularity plays a substantial role for what morphological properties are captured by a language model – and that this factor amplifies the effect of the language's share of the pre-training data.

It may seem unexpected that a higher degree of irregularity has a positive effect on probing accuracy. However, a possible explanation is that irregular forms are better memorized *because* they appear more often in the training data, given the correlation of irregularity with frequency (Wu et al., 2019). This could also be linked to the word predictability advantage of Zipfian distributions that has been shown to aid word segmentation in humans (Lavi-Rotbain and Arnon, 2022) – which we deem an interesting direction for future work.

Limitations and ethical considerations can be found in Appendices A and B, respectively.

# 7 Conclusion

We have analyzed the effects of tokenization, training data proportions, and linguistic factors on morphological knowledge encoded in the parameters of pre-trained multilingual language models. Through analyzing 17 languages and 4 morphological tasks, we have shown that the morphological knowledge encoded in multilingual language models differs across languages, despite the global average scores being similar. Beyond differences across languages, we also found differences across tasks, showing that tense and case are particularly hard to find in the representations of character-based language models. To further understand what exactly influences those difference, we have analyzed the effect of morphological complexity in relation to the language's proportion in the language model's pre-training data. We found that the degree of irregularity plays a significant role and amplifies the effect of training data, suggesting that more irregular languages benefit from a having a higher share in the data mix used for pre-training.

# References

Farrell Ackerman and Robert Malouf. 2013. Morphological organization: The low conditional entropy conjecture. *Language*, 89(3):429–464.

Judit Acs, Endre Hamerlik, Roy Schwartz, Noah A Smith, and Andras Kornai. 2023. Morphosyntactic probing of multilingual bert models. *Natural Language Engineering*, pages 1–40.

Judit Ács, Ákos Kádár, and Andras Kornai. 2021. Subword pooling makes a difference. In *Proceedings of the 16th Conference of the European Chapter of the Association for Computational Linguistics: Main Volume*, pages 2284–2295, Online. Association for Computational Linguistics.

Yossi Adi, Einat Kermany, Yonatan Belinkov, Ofer Lavi, and Yoav Goldberg. 2017. Fine-grained analysis of sentence embeddings using auxiliary prediction tasks. In *International Conference on Learning Representations*.

Mehdi Ali, Michael Fromm, Klaudia Thellmann, Richard Rutmann, Max Lübbering, Johannes Leveling, Katrin Klug, Jan Ebert, Niclas Doll, Jasper Buschhoff, Charvi Jain, Alexander Weber, Lena Jurkschat, Hammam Abdelwahab, Chelsea John, Pedro Ortiz Suarez, Malte Ostendorff, Samuel Weinbach, Rafet Sifa, Stefan Kesselheim, and Nicolas Flores-Herr. 2024. Tokenizer choice for LLM training: Negligible or crucial? In *Findings of the Association for Computational Linguistics: NAACL 2024*, pages 3907–3924, Mexico City, Mexico. Association for Computational Linguistics.

Dang Anh, Limor Raviv, and Lukas Galke. 2024. Morphology matters: Probing the cross-linguistic morphological generalization abilities of large language models through a Wug test. In *Proceedings of the Workshop on Cognitive Modeling and Computational Linguistics*, pages 177–188, Bangkok, Thailand. Association for Computational Linguistics.

Yonatan Belinkov. 2022. Probing classifiers: Promises, shortcomings, and advances. *Computational Linguistics*, 48(1):207–219.

Yonatan Belinkov, Nadir Durrani, Fahim Dalvi, Hassan Sajjad, and James Glass. 2017a. What do neural machine translation models learn about morphology? In *Proceedings of the 55th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 861–872, Vancouver, Canada. Association for Computational Linguistics.

Yonatan Belinkov, Nadir Durrani, Fahim Dalvi, Hassan Sajjad, and James Glass. 2020. On the linguistic representational power of neural machine translation models. *Computational Linguistics*, 46(1):1–52.

Yonatan Belinkov, Lluís Màrquez, Hassan Sajjad, Nadir Durrani, Fahim Dalvi, and James Glass. 2017b. Evaluating layers of representation in neural machine translation on part-of-speech and semantic tagging tasks. In *Proceedings of the Eighth International Joint Conference on Natural Language Processing (Volume 1: Long Papers)*.

Christian Bentz, Annemarie Verkerk, Douwe Kiela, Felix Hill, and Paula Buttery. 2015. Adaptive communication: Languages with more non-native speakers tend to have fewer word forms. *PloS one*, 10(6):e0128254.

Arianna Bisazza and Clara Tump. 2018. The lazy encoder: A fine-grained analysis of the role of morphology in neural machine translation. In *Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing*, pages 2871–2876, Brussels, Belgium. Association for Computational Linguistics.

Leonard Bloomfield. 1933. *Language*. H. Holt.

Piotr Bojanowski, Edouard Grave, Armand Joulin, and Tomas Mikolov. 2017. Enriching word vectors with subword information. *Transactions of the Association for Computational Linguistics*, 5:135–146.

Tom Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, et al. 2020. Language models are few-shot learners. *Advances in neural information processing systems*, 33:1877–1901.

Sébastien Bubeck, Varun Chandrasekaran, Ronen Eldan, Johannes Gehrke, Eric Horvitz, Ece Kamar, Peter Lee, Yin Tat Lee, Yuanzhi Li, Scott Lundberg, et al. 2023. Sparks of artificial general intelligence: Early experiments with GPT-4. *arXiv preprint arXiv:2303.12712*.

Junyoung Chung, Kyunghyun Cho, and Yoshua Bengio. 2016. A character-level decoder without explicit segmentation for neural machine translation. In *Proceedings of the 54th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 1693–1703, Berlin, Germany. Association for Computational Linguistics.

Jonathan H. Clark, Dan Garrette, Iulia Turc, and John Wieting. 2022. Canine: Pre-training an efficient tokenization-free encoder for language representation. *Transactions of the Association for Computational Linguistics*, 10:73–91.

Alexis Conneau, Kartikay Khandelwal, Naman Goyal, Vishrav Chaudhary, Guillaume Wenzek, Francisco Guzmán, Edouard Grave, Myle Ott, Luke Zettlemoyer, and Veselin Stoyanov. 2020. Unsupervised cross-lingual representation learning at scale. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 8440–8451, Online. Association for Computational Linguistics.

Alexis Conneau, German Kruszewski, Guillaume Lample, Loïc Barrault, and Marco Baroni. 2018. What you can cram into a single $&!#* vector: Probing

sentence embeddings for linguistic properties. In *Proceedings of the 56th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 2126–2136, Melbourne, Australia. Association for Computational Linguistics.

Ryan Cotterell, Sabrina J. Mielke, Jason Eisner, and Brian Roark. 2018. Are all languages equally hard to language-model? In *Proceedings of the 2018 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 2 (Short Papers)*, pages 536–541, New Orleans, Louisiana. Association for Computational Linguistics.

Robert M DeKeyser. 2005. What makes learning second-language grammar difficult? A review of issues. *Language learning*, 55.

Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. 2019. BERT: Pre-training of deep bidirectional transformers for language understanding. pages 4171–4186.

Matthew S Dryer and Martin Haspelmath. 2013. WALS Online (v2020. 3). *Zenodo*, 7385533.

Lukas Edman, Gabriele Sarti, Antonio Toral, Gertjan van Noord, and Arianna Bisazza. 2024. Are character-level translations worth the wait? Comparing ByT5 and mT5 for machine translation. *Transactions of the Association for Computational Linguistics*, 12:392–410.

Daniel Edmiston. 2020. A systematic analysis of morphological content in bert models for multiple languages. *arXiv preprint arXiv:2004.03032*.

Allyson Ettinger, Ahmed Elgohary, Colin Phillips, and Philip Resnik. 2018. Assessing composition in sentence vector representations. In *Proceedings of the 27th International Conference on Computational Linguistics*, pages 1790–1801, Santa Fe, New Mexico, USA. Association for Computational Linguistics.

Nicholas Evans and Stephen C Levinson. 2009. The myth of language universals: Language diversity and its importance for cognitive science. *Behavioral and brain sciences*, 32(5):429–448.

William Fleshman and Benjamin Van Durme. 2023. Toucan: Token-aware character level language modeling. *arXiv preprint arXiv:2311.08620*.

Lukas Galke, Yoav Ram, and Limor Raviv. 2024. Deep neural networks and humans both benefit from compositional language structure. *Nature Communications*, 15(10816).

Yingqiang Gao, Nikola I. Nikolov, Yuhuang Hu, and Richard H.R. Hahnloser. 2020. Character-level translation with self-attention. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 1591–1604, Online. Association for Computational Linguistics.

Daniela Gerz, Ivan Vulić, Edoardo Ponti, Jason Naradowsky, Roi Reichart, and Anna Korhonen. 2018. Language modeling for morphologically rich languages: Character-aware modeling for word-level prediction. *Transactions of the Association for Computational Linguistics*, 6:451–465.

Omer Goldman and Reut Tsarfaty. 2022. Morphology without borders: Clause-level morphology. *Transactions of the Association for Computational Linguistics*, 10:1455–1472.

John Hewitt, Kawin Ethayarajh, Percy Liang, and Christopher Manning. 2021. Conditional probing: Measuring usable information beyond a baseline. In *Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing*, pages 1626–1639, Online and Punta Cana, Dominican Republic. Association for Computational Linguistics.

John Hewitt and Percy Liang. 2019. Designing and interpreting probes with control tasks. In *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)*, pages 2733–2743, Hong Kong, China. Association for Computational Linguistics.

Valentin Hofmann, Janet Pierrehumbert, and Hinrich Schütze. 2021. Superbizarre is not superb: Derivational morphology improves BERT's interpretation of complex words. In *Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing (Volume 1: Long Papers)*, pages 3594–3608, Online. Association for Computational Linguistics.

Dieuwke Hupkes, Mario Giulianelli, Verna Dankers, Mikel Artetxe, Yanai Elazar, Tiago Pimentel, Christos Christodoulopoulos, Karim Lasri, Naomi Saphra, Arabella Sinclair, Dennis Ulmer, Florian Schottmann, Khuyagbaatar Batsuren, Kaiser Sun, Koustuv Sinha, Leila Khalatbari, Maria Ryskina, Rita Frieske, Ryan Cotterell, and Zhijing Jin. 2023. A taxonomy and review of generalization research in NLP. *Nature Machine Intelligence*, 5(10):1161–1174.

Dieuwke Hupkes, Sara Veldhoen, and Willem Zuidema. 2018. Visualisation and'diagnostic classifiers' reveal how recurrent and recursive neural networks process hierarchical structure. *Journal of Artificial Intelligence Research*, 61:907–926.

Vera Kempe and Patricia J Brooks. 2008. Second language learning of complex inflectional systems. *Language Learning*, 58(4):703–746.

Yoon Kim, Yacine Jernite, David Sontag, and Alexander Rush. 2016. Character-aware neural language models. In *Proceedings of the AAAI conference on artificial intelligence*, volume 30.

Taku Kudo. 2018. Subword regularization: Improving neural network translation models with multiple

[subword candidates](). In *Proceedings of the 56th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 66–75, Melbourne, Australia. Association for Computational Linguistics.

Taku Kudo and John Richardson. 2018. Sentencepiece: A simple and language independent subword tokenizer and detokenizer for neural text processing. In *Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing: System Demonstrations*, pages 66–71.

Guillaume Lample and Alexis Conneau. 2019. Crosslingual language model pretraining. *arXiv preprint arXiv:1901.07291*.

Ori Lavi-Rotbain and Inbal Arnon. 2022. The learnability consequences of zipfian distributions in language. *Cognition*, 223:105038.

Teven Le Scao, Angela Fan, Christopher Akiki, Ellie Pavlick, Suzana Ilić, Daniel Hesslow, Roman Castagné, Alexandra Sasha Luccioni, François Yvon, Matthias Gallé, et al. 2023. Bloom: A 176b-parameter open-access multilingual language model.

Jason Lee, Kyunghyun Cho, and Thomas Hofmann. 2017. Fully character-level neural machine translation without explicit segmentation. *Transactions of the Association for Computational Linguistics*, 5:365–378.

Yongjie Lin, Yi Chern Tan, and Robert Frank. 2019. Open sesame: Getting inside bert's linguistic knowledge. In *Proceedings of the 2019 ACL Workshop BlackboxNLP: Analyzing and Interpreting Neural Networks for NLP*, pages 241–253.

Nelson F Liu, Matt Gardner, Yonatan Belinkov, Matthew E Peters, and Noah A Smith. 2019. Linguistic knowledge and transferability of contextual representations. In *Proceedings of NAACL-HLT*, pages 1073–1094.

Gary Lupyan and Rick Dale. 2010. [Language Structure Is Partly Determined by Social Structure](). *PLoS ONE*, 5(1):e8559.

Andreas Madsen, Siva Reddy, and Sarath Chandar. 2022. Post-hoc interpretability for neural nlp: A survey. *ACM Computing Surveys*, 55(8):1–42.

Christopher D Manning, Kevin Clark, John Hewitt, Urvashi Khandelwal, and Omer Levy. 2020. Emergent linguistic structure in artificial neural networks trained by self-supervision. *Proceedings of the National Academy of Sciences*, 117(48):30046–30054.

Sabrina J Mielke, Ryan Cotterell, Kyle Gorman, Brian Roark, and Jason Eisner. 2019. What kind of language is hard to language-model? In *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*, pages 4975–4989.

Hyunji Hayley Park, Katherine J Zhang, Coleman Haley, Kenneth Steimel, Han Liu, and Lane Schwartz. 2021. Morphology matters: A multilingual language modeling analysis. *Transactions of the Association for Computational Linguistics*, 9:261–276.

Matthew E Peters, Mark Neumann, Luke Zettlemoyer, and Wen-tau Yih. 2018. Dissecting contextual word embeddings: Architecture and representation. In *Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing*, pages 1499–1509.

Peng Qian, Xipeng Qiu, and Xuan-Jing Huang. 2016. Investigating language universal and specific properties in word embeddings. In *Proceedings of the 54th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 1478–1488.

Colin Raffel, Noam Shazeer, Adam Roberts, Katherine Lee, Sharan Narang, Michael Matena, Yanqi Zhou, Wei Li, and Peter J Liu. 2020. Exploring the limits of transfer learning with a unified text-to-text transformer. *Journal of machine learning research*, 21(140):1–67.

Limor Raviv, Marianne de Heer Kloots, and Antje Meyer. 2021. What makes a language easy to learn? a preregistered study on how systematic structure and community size affect language learnability. *Cognition*, 210:104620.

Limor Raviv, Louise R Peckre, and Cedric Boeckx. 2022. What is simple is actually quite complex: A critical note on terminology in the domain of language and communication. *Journal of Comparative Psychology*, 136(4):215.

Anna Rogers, Olga Kovaleva, and Anna Rumshisky. 2021. A primer in bertology: What we know about how bert works. *Transactions of the Association for Computational Linguistics*, 8:842–866.

Sebastian Ruder, Ivan Vulić, and Anders Søgaard. 2019. A survey of cross-lingual word embedding models. *Journal of Artificial Intelligence Research*, 65:569–631.

Benoît Sagot. 2013. Comparing complexity measures. In *Computational approaches to morphological complexity*.

Rico Sennrich, Barry Haddow, and Alexandra Birch. 2016. Neural machine translation of rare words with subword units. In *Proceedings of the 54th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 1715–1725.

Anders Søgaard, Anders Johannsen, Barbara Plank, Dirk Hovy, and Hector Martínez Alonso. 2014. What's in a p-value in nlp? In *Proceedings of the eighteenth conference on computational natural language learning*, pages 1–10.

Ian Tenney, Dipanjan Das, and Ellie Pavlick. 2019. Bert rediscovers the classical nlp pipeline. In *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*, pages 4593–4601.

Olia Toporkov and Rodrigo Agerri. 2024. On the role of morphological information for contextual lemmatization. *Computational Linguistics*, pages 1–35.

Hugo Touvron, Thibaut Lavril, Gautier Izacard, Xavier Martinet, Marie-Anne Lachaux, Timothée Lacroix, Baptiste Rozière, Naman Goyal, Eric Hambro, Faisal Azhar, et al. 2023. Llama: Open and efficient foundation language models. *arXiv preprint arXiv:2302.13971*.

Ekaterina Vylomova, Trevor Cohn, Xuanli He, and Gholamreza Haffari. 2017. Word representation models for morphologically rich languages in neural machine translation. In *Proceedings of the First Workshop on Subword and Character Level Models in NLP*, pages 103–108.

Alex Warstadt, Yian Zhang, Xiaocheng Li, Haokun Liu, and Samuel Bowman. 2020. Learning which features matter: Roberta acquires a preference for linguistic generalizations (eventually). In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pages 217–235.

Jason Wei, Yi Tay, Rishi Bommasani, Colin Raffel, Barret Zoph, Sebastian Borgeaud, Dani Yogatama, Maarten Bosma, Denny Zhou, Donald Metzler, Ed H. Chi, Tatsunori Hashimoto, Oriol Vinyals, Percy Liang, Jeff Dean, and William Fedus. 2022. Emergent abilities of large language models. *Transactions on Machine Learning Research*. Survey Certification.

Shijie Wu, Ryan Cotterell, and Timothy O'Donnell. 2019. Morphological irregularity correlates with frequency. In *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*, pages 5117–5126.

Shijie Wu and Mark Dredze. 2020. Are all languages created equal in multilingual BERT? In *Proceedings of the 5th Workshop on Representation Learning for NLP*, pages 120–130.

Linting Xue, Aditya Barua, Noah Constant, Rami Al-Rfou, Sharan Narang, Mihir Kale, Adam Roberts, and Colin Raffel. 2022. ByT5: Towards a token-free future with pre-trained byte-to-byte models. *Transactions of the Association for Computational Linguistics*, 10:291–306.

Linting Xue, Noah Constant, Adam Roberts, Mihir Kale, Rami Al-Rfou, Aditya Siddhant, Aditya Barua, and Colin Raffel. 2021. mT5: A massively multilingual pre-trained text-to-text transformer. In *Proceedings of the 2021 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pages 483–498, Online. Association for Computational Linguistics.

Kelly Zhang and Samuel Bowman. 2018. Language modeling teaches you more than translation does: Lessons learned through auxiliary syntactic task analysis. In *Proceedings of the 2018 EMNLP Workshop BlackboxNLP: Analyzing and Interpreting Neural Networks for NLP*, pages 359–361.

## A    Limitations

Several limitations should be taken into account: First, our experiments do not take orthographic transparency into account, as all models and complexity measures are based on written language. Second, we ran each experiment only once due to the high volume of experiments. Third, while we aimed to cover as many typologically different languages as possible, the dataset we use supports mostly Indo-European languages, such that 12 out of our 17 considered languages are Indo-European. Moreover, the current study only focuses on inflectional morphology due to the lack of datasets for probing derivational morphology. Lastly, p-values should be treated with care when the sample size (here: 17,200) is large (Søgaard et al., 2014).

## B    Ethical Considerations

We emphasize that morphological complexity of languages bears no implication on their quality – having more complexity does not make one language better than another (see Raviv et al., 2022).

## C    Effect of Probing Architecture

Previous studies on probing linguistic features have had a debate over which type of probe is sufficient to extract the relevant knowledge, but insufficient to learn the knowledge itself (Belinkov, 2022). Here, we also compare the accuracy scores of MLP probes and linear probes for mT5 and ByT5, as shown in Table 4.

| Model | Probing Architecture | |
| --- | --- | --- |
| | **MLP** | **Linear** |
| mT5-base | 82.57 | 80.37 |
| ByT5-base | 82.86 | 80.61 |

Table 4: Probing accuracy of mT5 and ByT5 when using MLP and linear classifiers, averaged over languages and tasks

We observed no considerable differences in accuracy scores across types of probes. For both types of probes, the mean accuracy scores are all around

80%. This pattern is most inline with Acs et al. (2023), suggesting that linear classifiers are as effective as non-linear ones in extracting morphological knowledge of multilingual LLMs. Moreover, the observation that linear probes perform equally well as MLPs implies that morphology is a relatively simple feature that can be learned early and straightforwardly by the models.

## D Effect of Pooling Methods

We compare the probing accuracy when using these two pooling methods, namely the **average** and **last** method. Table 5 shows the results of when using these two methods for mT5 and ByT5.

| Model | Subword Pooling Method | |
|---|---|---|
| | **last** | **average** |
| mT5-base | 82.57 | 75.28 |
| ByT5-base | 82.86 | 76.40 |

Table 5: Probing accuracy of mT5 and ByT5 when using different subword pooling methods. The results were averaged over languages and tasks and the reported probes are MLPs

Comparing the accuracy scores of the two pooling methods, it can be seen that the *last* method achieved considerably higher accuracy scores than the *average* method. The difference is approximately 6-7 points. This also holds for both models. It seems that the representational information and/or the morphological content of a word is mostly encoded in its last token. Previous comparisons have also reported similar results (Acs et al., 2023; Ács et al., 2021; Belinkov et al., 2020).

## E Details of the Considered Languages

### E.1 Morphological Properties

Table 6 shows the morphological properties of the considered languages.

### E.2 Proportion of mT5/ByT5's training data and morphological complexity of the language

Table 7 shows the language's proportion of training data and their morphological complexity scores: TTR and Irregularity.

| Language | Tr. Data | TTR | Irregularity |
|---|---|---|---|
| English | 5.67% | -0.460 | -5.94 |
| German | 3.05% | -0.010 | -6.28 |
| Dutch | 1.98% | -0.390 | -6.68 |
| French | 2.89% | -0.340 | -4.16 |
| Romanian | 1.58% | -0.420 | -3.40 |
| Spanish | 3.09% | 0.001 | -8.81 |
| Portuguese | 2.36% | 0.038 | -9.11 |
| Turkish | 1.93% | 1.550 | -5.96 |
| Czech | 1.72% | 0.430 | -5.63 |
| Russian | 3.71% | 0.870 | -7.74 |
| Hebrew | 1.06% | 2.020 | -1.78 |
| Arabic | 1.66% | 1.630 | -0.06 |
| Hindi | 1.21% | -0.300 | -2.10 |
| Estonian | 0.89% | 1.760 | -2.79 |
| Latvian | 0.87% | 0.770 | -7.90 |
| Urdu | 0.61% | -0.450 | 9.20 |
| Basque | 0.57% | 1.310 | 19.86 |

Table 7: Percentages of training data of mT5 and Byt5 from Xue et al. (2021), irregularity scores from Wu et al. (2019), and TTR scores from Bentz et al. (2015) for each investigated language. For irregularity, higher scores mean being more morphologically irregular. In contrast, higher TTRs mean higher complexity.

## F Considered Probing Tasks

Here, we provides an overview the morphological properties that we investigate in the study, namely number, tense, case, and gender, and how they may vary between languages. (Acs et al., 2023)

**Number** In many languages, especially inflected languages, nouns are marked as either singular or plural (Bloomfield, 1933). An exception is Latvian, which includes singular, plural, and partitive nouns. Plurality is usually expressed by adding certain endings to the nouns, and sometimes include changing their vowels. These endings are determined in different ways across languages. For instance, in some Indo-European languages such as Spanish, the plural form of nouns is affected by their gender. In this task, the LLMs have to predict whether the target word is a plural or singular noun.

**Tense** Most languages mark tenses (Bloomfield, 1933). In some languages, tenses are indicated by inflecting verbs. In other languages, for example, Estonian, adjectives can also express tense. In certain languages, tense can interact with other morphological features, namely mood and aspect. Inflection patterns for tense are usually dependent

on the ending, conjugation pattern of the verbs, and whether they are regular or irregular. In Spanish and French, it is also dependent on the subject pronouns. In Hindi, verbs indicating tense must agree with gender and number of the subject.

**Case**   A case system is a grammatical category used in many languages to mark the relationship between a noun or pronoun and other words in a sentence. Case marking is typically indicated through inflection. The number of cases varies across languages. Cases are often marked with inflection. In some languages, case often affects how articles, pronouns, and adjectives should be inflected. Previous probing studies show that case is often one of the most challenging morphological categories to be learned by language models (Edmiston, 2020; Bisazza and Tump, 2018; Acs et al., 2023).

**Gender**   Most Indo-European languages mark genders in nouns and often require agreement with in other part-of-speech in the sentence, such as verbs and adjectives Bloomfield (1933). Gender systems exhibit substantial diversity in their number of genders, assignment rules. Some languages (e.g., Basque) do not have a gender system. Romance languages have a binary gender system. On the other hand, Germanic and Slavic languages often have more than two genders. For example, Dutch nouns are either common or neuter while German nouns can be masculine, feminine, or neuter. Spanish masculine nouns end in *"-o"* while feminine nouns end in *"a"*. There is also a certain degree of irregularity.

## G   Extended Results

Table 8 shows the full breakdown of probing accuracy per task and per language.

## H   Detailed Results of the Statistical Analysis

We provide the formulation and results of the statistical models. We define our model as the combination of two interaction terms between `irregularity` and `TD` and between `TTR` and `TD`. The syntax is as follows. The results can be found in Table 9.

```
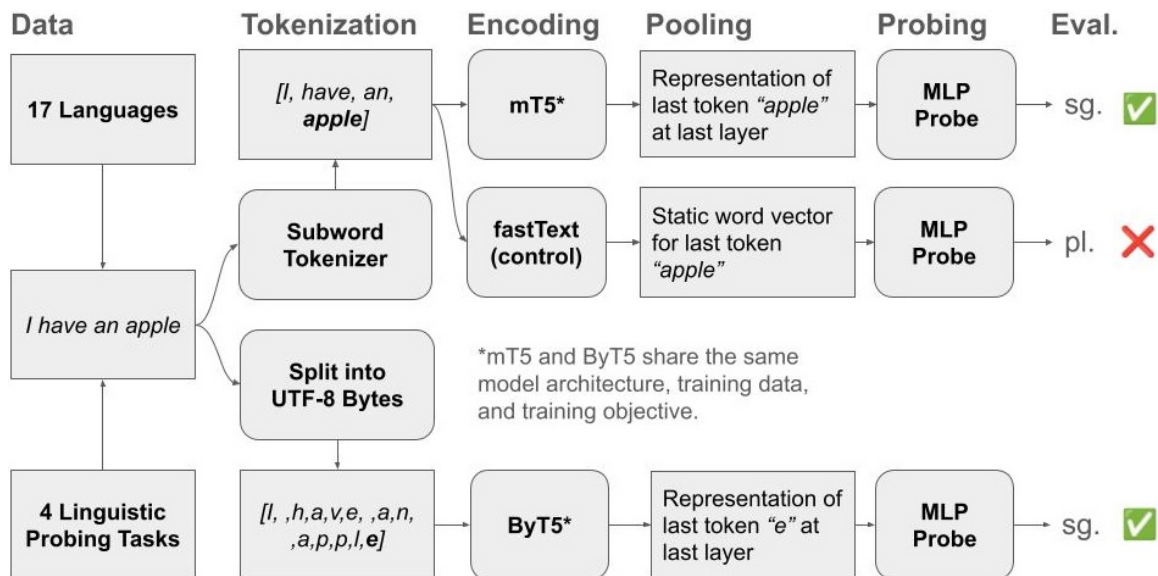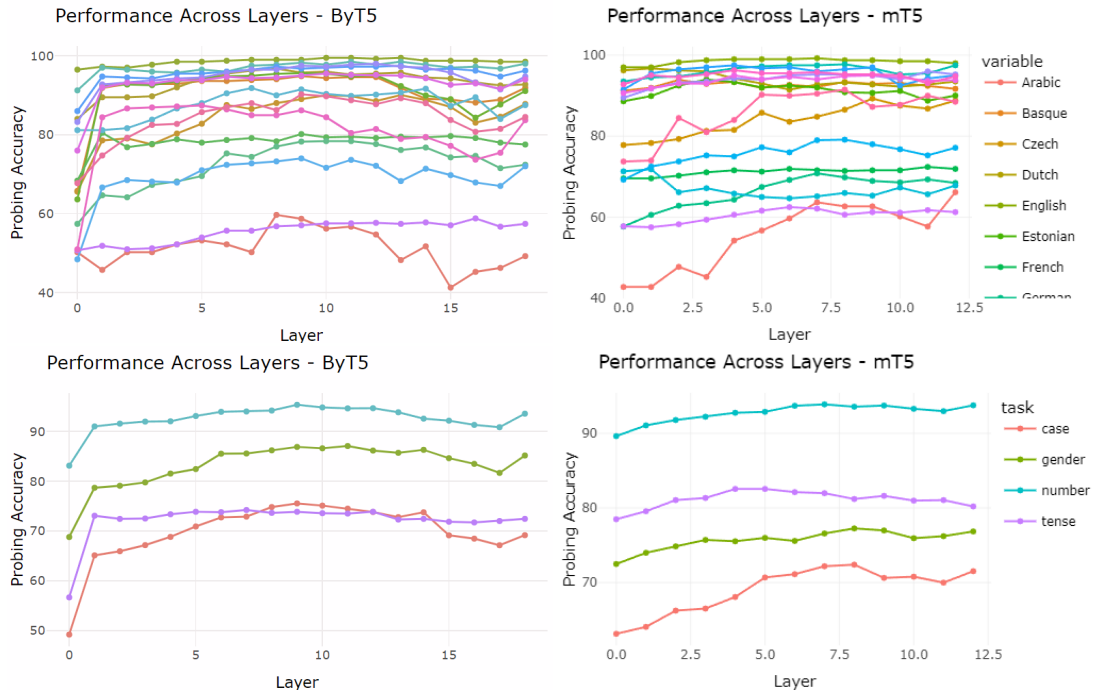accuracy ~ irregularity*TD + TTR*TD +
(1|language) + (1|task)
```

| **Fixed Effects** | | | | |
|---|---|---|---|---|
| Variable | Estimate | SE | *t*-value | *p*-value |
| (Intercept) | 2.71845 | 0.61012 | 4.456 | <.001 *** |
| Training data (TD) | 2.11411 | 0.55564 | 3.805 | **<.001** *** |
| Irregularity (I) | 2.83197 | 0.42242 | 6.704 | **<.001** *** |
| TTR (T) | -0.62728 | 0.57127 | -1.098 | 0.272179 |
| TD*I | 2.03621 | 0.31730 | 6.417 | **<.001** *** |
| TD*T | -0.62728 | 0.57127 | -1.098 | 0.272179 |
| **Random Effect** | | | | |
| Group | Name | Variance | Std.Dev. | |
| language | (Intercept) | 3.0920 | 1.7584 | |
| task | (Intercept) | 0.3571 | 0.5976 | |

Table 9: Results of linear mixed effect regression with *probing accuracy* of mT5 as outcome variable and *language* and *task* as random effects. Fixed effects are irregularity (I) and TTR (T), training data (TD) and their two-way interactions.

| Language | Genus | Family | POS | Number | Tense | Case | Gender |
|----------|-------|--------|-----|--------|-------|------|--------|
| English | Germanic | Indo-European | N | 2 | – | – | – |
| English | Germanic | Indo-European | V | – | 2 | – | – |
| German | Germanic | Indo-European | N | 2 | – | 4 | 3 |
| German | Germanic | Indo-European | V | – | 2 | – | – |
| Dutch | Germanic | Indo-European | N | 2 | – | – | 2 |
| French | Romance | Indo-European | N | 2 | – | – | 2 |
| French | Romance | Indo-European | V | – | 4 | – | – |
| Spanish | Romance | Indo-European | N | 2 | – | – | 2 |
| Spanish | Romance | Indo-European | V | – | 4 | – | – |
| Portuguese | Romance | Indo-European | N | 2 | – | – | 2 |
| Romanian | Romance | Indo-European | N | 2 | – | – | 2 |
| Turkish | Turkic | Altaic | N | 2 | – | 7 | – |
| Russian | Slavic | Indo-European | N | 2 | – | 6 | 3 |
| Russian | Slavic | Indo-European | V | – | 3 | – | – |
| Czech | Slavic | Indo-European | N | 2 | – | – | 3 |
| Hebrew | Semitic | Afro-Asiatic | N | 2 | – | – | 2 |
| Hindi | Indic | Indo-European | N | 2 | – | 2 | 2 |
| Urdu | Indic | Indo-European | N | 2 | – | 2 | – |
| Urdu | Indic | Indo-European | V | 2 | – | – | – |
| Basque | Basque | Basque | N | 2 | – | 11 | – |
| Estonian | Finnic | Uralic | N | 2 | – | 18 | – |
| Estonian | Finnic | Uralic | V | – | – | – | – |
| Latvian | Baltic | Indo-European | N | 3 | – | 5 | 3 |
| Latvian | Baltic | Indo-European | V | – | 3 | – | – |
| Arabic | Semitic | Afro-Asiatic | N | – | – | 2 | – |

Table 6: List of studied languages, their genera and families, along with the number of possible classes within the studied dataset per morphological properties of 17 investigated languages (N = noun; V = verb)

| No. | Language | Task | mT5 | ByT5 | fastText |
|---|---|---|---|---|---|
| 1 | Arabic | case | 66.17 | 49.25 | 37.81 |
| 2 | Basque | case | 91.39 | 92.55 | 17.22 |
| 3 | Basque | number | 92.00 | 89.42 | 87.99 |
| 4 | Czech | gender | 81.09 | 78.87 | 72.13 |
| 5 | Czech | number | 96.50 | 90.05 | 85.50 |
| 6 | Dutch | gender | 90.00 | 88.39 | 72.50 |
| 7 | Dutch | number | 97.00 | 98.24 | 92.00 |
| 8 | English | number | 97.50 | 98.55 | 98.50 |
| 9 | English | tense | 98.50 | 98.53 | 97.00 |
| 10 | Estonian | case | 88.10 | 86.84 | 62.85 |
| 11 | Estonian | number | 91.50 | 92.47 | 89.49 |
| 12 | Estonian | tense | 90.50 | 93.84 | 96.49 |
| 13 | French | gender | 95.00 | 90.39 | 92.00 |
| 14 | French | number | 99.50 | 98.24 | 95.49 |
| 15 | French | tense | 21.29 | 46.30 | 91.08 |
| 16 | German | case | 65.00 | 40.42 | 28.00 |
| 17 | German | gender | 29.85 | 74.42 | 76.00 |
| 18 | German | number | 92.00 | 89.00 | 84.50 |
| 19 | German | tense | 87.00 | 85.95 | 87.00 |
| 20 | Hebrew | gender | 95.50 | 95.05 | 89.49 |
| 21 | Hebrew | number | 99.50 | 98.82 | 95.49 |
| 22 | Hindi | case | 13.00 | 81.58 | 63.49 |
| 23 | Hindi | gender | 95.50 | 90.74 | 49.00 |
| 24 | Hindi | number | 95.00 | 90.95 | 61.57 |
| 25 | Latvian | case | 89.50 | 98.74 | 84.50 |
| 26 | Latvian | gender | 68.00 | 32.61 | 64.49 |
| 27 | Latvian | number | 68.50 | 70.39 | 63.49 |
| 28 | Latvian | tense | 82.59 | 84.50 | 82.58 |
| 29 | Portuguese | gender | 95.00 | 93.66 | 97.00 |
| 30 | Portuguese | number | 95.00 | 97.55 | 97.50 |
| 31 | Romanian | gender | 93.50 | 93.68 | 91.00 |
| 32 | Romanian | number | 97.00 | 95.84 | 93.50 |
| 33 | Russian | case | 48.04 | 36.02 | 82.55 |
| 34 | Russian | gender | 8.46 | 82.98 | 50.74 |
| 35 | Russian | number | 96.00 | 93.39 | 91.50 |
| 36 | Russian | tense | 92.54 | 9.59 | 92.03 |
| 37 | Spanish | gender | 93.50 | 94.66 | 97.00 |
| 38 | Spanish | number | 99.00 | 97.84 | 97.50 |
| 39 | Spanish | tense | 89.00 | 86.47 | 31.00 |
| 40 | Turkish | case | 95.57 | 72.88 | 61.57 |
| 41 | Turkish | number | 94.00 | 89.32 | 94.99 |
| 42 | Urdu | case | 87.00 | 77.37 | 61.50 |
| 43 | Urdu | number | 90.00 | 90.87 | 81.49 |

Table 8: Accuracy scores of each task in each language from mT5, ByT5, and fastText

# Dora explores Clinically Relevant Information in EHRs using NER

**Martin Sundahl Laursen**
Department of Clinical Biochemistry
Odense University Hospital
`martin.sundahl.laursen@rsyd.dk`

**Lina Elkjær Pedersen**
Department of Clinical Biochemistry
Odense University Hospital

**Josefine Bak H Adelhelm**
Department of Clinical Biochemistry
Odense University Hospital

**Rasmus Bank Lynggaard**
Department of Clinical Biochemistry
Odense University Hospital

**Pernille Just Vinholt**
Department of Clinical Biochemistry, Odense University Hospital
Department of Clinical Research, University of Southern Denmark

## Abstract

Retrieving relevant information from unstructured electronic health records is time-consuming and prone to error, reducing time available for direct patient care. We present Dora, a Danish clinical named entity recognition model that builds on prior work by Laursen et al. (2023a). Dora identifies six types of clinical entities to support medical information retrieval: diseases, symptoms/findings, diagnostics, treatments, anatomies, and results. The model achieves an exact boundary macro F1 score of 0.922 and overlap boundary score of 0.945. A prospective clinical utility evaluation shows that Dora reliably extracts relevant information for physicians. A bias analysis indicates slightly reduced performance on psychiatric notes, with minimal overall differences.

## 1 Introduction

Health care professionals, particularly medical doctors (MDs), need to retrieve information from electronic health records (EHRs) regarding diagnoses, symptoms, medications, treatments, etc. This process is time-consuming, carries the risk of overlooking important information, and ultimately reduces the time available for direct patient care (Laursen et al., 2023b). Furthermore, the health care data is in an unstructured format in the EHR. The EHR system may include a basic "find on page" function, which allows users to search for specific words or phrases within the visible text. However, this method is vulnerable to inaccuracies such as misspellings, abbreviations, and typographical errors.

Previous work has shown that natural language processing methods, particularly Named Entity Recognition (NER) models, can effectively identify clinical entities in EHR text (Jiang et al., 2011; Alsentzer et al., 2019).

Notable results in English clinical NER are Stanza (Qi et al., 2020; Zhang et al., 2021) with micro F1 0.881 and BioBERT (Lee et al., 2020) with micro F1 0.867 on identifying problems, tests, and treatments in the i2b2 dataset (Uzuner et al., 2011). In Scandinavian clinical NER, Laursen et al. (2023a) achieved an entity-level macro F1 of 0.601 for exact boundary matching and 0.682 for overlap boundaries when identifying diseases, symptoms (including abnormal findings), diagnostics, treatments, anatomies, and results in Danish EHR text. In Swedish, RoBERTa Large by AI Sweden achieved a token-level micro F1 of 0.779 when compared to eight other encoder models on identifying diagnoses, findings, body parts, and drugs in the Stockholm EPR Clinical Entity Corpus (Vakili et al., 2025; Skeppstedt et al., 2014).

Few medical machine learning studies, however, extend beyond reporting internal test set performance and do not assess real-world clinical impact and utility (Kelly et al., 2019; Ghassemi et al., 2020; Rajpurkar et al., 2022).

In this paper, we extend the Danish Clinical NER model by Laursen et al. (2023a), retraining it on an expanded and re-annotated dataset with updates to the annotation scheme, preprocessing, postprocessing, and model training.

We present a new Danish clinical NER model, Dora, that presents substantial improvements in model performance. We demonstrate the clinical utility in a prospective real-world evaluation and evaluate bias.

## 2 Methods

In this section, we first describe the data sources for the model's development and evaluation cohort. We then outline the model's development

and present the different evaluations conducted to assess the model.

## 2.1 Data Sources

We used data from two different EHR systems of Odense University Hospital in the Region of Southern Denmark for development and evaluation. The COSMIC cohort consisted of EHRs from the COSMIC system (Cambio, CGI, Denmark) from November 2015 to September 2020. The EPJ cohort consisted of all EHRs from the EPJ SYD (Systematic, Denmark) system from February 2022 to November 2023.

## 2.2 Model Development

The development of the first iteration of the NER model was previously described by Laursen et al. (2023a). Here, we focus on refinements made to the annotation scheme, dataset, system architecture, and model development.

### 2.2.1 Annotation

We built on the clinical event annotation scheme proposed by Laursen et al. (2023a), with one key revision to improve usability for healthcare professionals: the Symptom entity, which used to include symptoms and pathological findings, now includes symptoms and all clinical findings—either normal or pathological.

The dataset from Laursen et al. (2023a) was re-annotated by a MD to reflect the revised scheme and iteratively extended with paragraphs from the COSMIC and EPJ cohorts using active learning and a locally developed annotation tool. Targeted data augmentation was applied to address specific errors.

### 2.2.2 Dataset

Our dataset contained 158,839 total entities, almost triple the size of the original dataset, split into training, validation and test sets. Splits were stratified to maintain a balanced distribution across entity labels, see Table 1.

### 2.2.3 System Architecture

We adapted the Princeton University Relation Extraction system (PURE) (Zhong and Chen, 2021), using code with minor modifications from Laursen et al. (2023a).

PURE classifies entities from constructed span-embeddings, which is the concatenation of the contextual embeddings of the start and end tokens,

along with a learned span-width embedding (Zhong and Chen, 2021). For full architectural details, we refer to the original work. Our modifications applied to the implementation by Laursen et al. (2023a) include:

- **Preprocessing:** Lowercasing, removing non-printable/control characters, converting HTML entities to Unicode, and mapping uncommon accented or special characters to standard equivalents.

- **Postprocessing:** After prediction, overlapping spans with the same label are merged. When overlaps have different labels, a voting mechanism selects the most likely label.

### 2.2.4 Development

We followed Laursen et al. (2023a) in extracting contextual embeddings using a Danish clinical ELECTRA encoder (Pedersen et al., 2022; Clark et al., 2020). Spans ranged from 1–10 tokens. Each of the start, end, and width components had size 256.

We trained using AdamW (Loshchilov and Hutter, 2019) (weight decay 0.001, batch size 32), early stopping (patience 6), and learning rate scheduling (patience 3, factor 0.2). No class weighting was applied.

Optimal learning rates (search space in parenthesis) were 5e-5 (5e-6–7.5e-5) for the encoder and 5e-4 (5e-5–7.5e-4) for the classifier. The optimal span classifier configuration was one (0–2) 1024-unit (256–1024; plateaued at 1024) hidden layer with ReLU activation and 0.3 dropout.

Model selection was based on the span-level macro F1 score on the validation set, excluding the negative class. We report the best model's entity-level recall, precision, and F1 score on the test set, using exact and overlapping boundary matching (Chinchor and Sundheim, 1993). We report a confusion matrix based on overlap matching, which better reflects clinical utility due to the often ambiguous boundaries of clinical entities.

## 2.3 Clinical Evaluation

The aim of the clinical evaluation was to assess the model's clinical utility and potential bias on an evaluation cohort stratified by gender (male/female), age group (child/adult/senior), and 15 diagnoses. To ensure diverse diagnoses, two MDs selected five diagnoses within each medical area; medical, psychiatry, surgical, from The Danish Med-

|  | Train (% of row total) | Validation (% of row total) | Test (% of row total) | TOTAL (% of column total) |
|---|---|---|---|---|
| **Paragraphs** | 18,001 (80%) | 2,206 (10%) | 2,258 (10%) | 22,465 (100%) |
| | | **Clinical events** | | |
| **Disease** | 7,198 (81%) | 821 (9%) | 887 (10%) | 8,906 (6%) |
| **Symptom** | 37,692 (80%) | 4,467 (10%) | 4,808 (10%) | 46,967 (30%) |
| **Treatment** | 22,218 (80%) | 2,806 (10%) | 2,774 (10%) | 27,798 (18%) |
| **Diagnostic** | 21,631 (80%) | 2,782 (10%) | 2,654 (10%) | 27,067 (17%) |
| **Anatomy** | 25,444 (80%) | 3,104 (10%) | 3,234 (10%) | 31,782 (20%) |
| **Result** | 13,024 (80%) | 1,714 (11%) | 1,581 (10%) | 16,319 (10%) |
| **TOTAL** | 127,207 (80%) | 15,694 (10%) | 15,938 (10%) | 158,839 (100%) |

Table 1: Distribution of clinical event types in the training, validation, and test sets.

ical Classification System (SKS) (Danish Health Data Authority, n.d.), which is based on the International Classification of Diseases 10th revision (ICD-10) (World Health Organization, 2016). Diagnoses spanning all genders and age groups and with a high likelihood of mentions of varied clinical entities like symptoms, diagnostics, and treatments were chosen:

- **Medical:** asthma, diabetic ketoacidosis (type 1), epilepsy, pneumonia, rheumatoid arthritis

- **Psychiatry:** autism, depression, eating disorder, generalised anxiety disorder, suicide attempt/self-harm

- **Surgical:** appendicitis, hernia, ileus, epistaxis, tibia fracture

We then randomly sampled EHRs from the EPJ cohort that included either a ICD-10 code or a textual mention of one of these 15 diagnoses.

### 2.3.1 Clinical Utility Evaluation

To evaluate the model's clinical utility, a MD manually reviewed its output on the evaluation cohort. For each EHR, the model's extracted entities were shown in a spreadsheet containing one row per entity with its label and context window. The full EHR text was provided for reference. EHRs were included iteratively, seeking a stratified sample of three different EHRs for each combination of diagnosis, gender, and age (n=270).

The MD assessed whether the model output included at least one mention of: 1) the disease entity for the target diagnosis, 2) symptoms, 3) diagnostic tool, and 4) treatment relevant for that diagnosis. If the diagnosis was missing from predictions, the full EHR was reviewed to confirm its presence. If the diagnosis was absent, the EHR was not included. When any expected entity was missing, the full

EHR was checked to determine if the model had failed to extract it.

257 samples were included. The cohort consisted of 132 females and 125 males, including 85 children, 86 adults, and 86 seniors. Two groups were entirely absent: female children with depression and senior males with eating disorder.

We calculated the detection rate per entity label.

The 15 diagnoses and expected clinical findings for each entity are presented in Appendix B.

### 2.3.2 Bias analysis

We conducted a structured bias analysis across gender, age group, and medical area.

From the evaluation cohort, we sampled three random patient EHRs (>5 notes available) per combination (n=270). Each patient was represented by four random notes (>49 characters per note to avoid minimal or templated content) (n=1,080).

Model predictions were corrected by a MD to establish ground truth. Entity-level F1 scores were calculated per patient, with micro and macro averages across labels (Chinchor and Sundheim, 1993). We report summary statistics for entity counts by medical area.

To ensure robust metrics given the short text span per patient, we applied conservative rules when averaging to handle missing entities:

- **No ground truths, some predictions:** recall excluded; precision and F1 set to 0.

- **Some ground truths, no predictions:** precision excluded; recall and F1 set to 0.

- **No ground truths or predictions:** all metrics excluded.

We further bootstrapped with 9,999 resamples per individual variable (i.e., each gender, age group,

and medical area) to produce 95% confidence intervals (CIs) by entity label and micro and macro average, using these to assess systematic model bias (Steyerberg et al., 2001).

## 3 Results

This section presents the results of the evaluation of the model performance, clinical utility, and potential biases.

### 3.1 Test Set Performance

The model achieved F1 scores above 0.90 across all entity types and evaluations. Macro F1 was 0.922 with exact boundary matching and 0.945 with overlap. Ignoring labels, the detection macro F1 with overlap reached 0.962. Detailed results are shown in Table 2.

| | TEST SET | | | | | |
|---|---|---|---|---|---|---|
| | **Exact boundary** | | | **Overlap boundary** | | |
| | **F1** | **Prec** | **Recall** | **F1** | **Prec** | **Recall** |
| **Disease** | 0.914 | 0.921 | 0.906 | 0.936 | 0.939 | 0.932 |
| **Symptom** | 0.902 | 0.917 | 0.888 | 0.930 | 0.943 | 0.918 |
| **Treatment** | 0.926 | 0.932 | 0.920 | 0.953 | 0.957 | 0.949 |
| **Diagnostic** | 0.941 | 0.943 | 0.938 | 0.957 | 0.958 | 0.956 |
| **Anatomy** | 0.940 | 0.950 | 0.930 | 0.968 | 0.974 | 0.962 |
| **Result** | 0.907 | 0.910 | 0.905 | 0.930 | 0.931 | 0.929 |
| **Micro avg** | 0.922 | 0.930 | 0.913 | 0.946 | 0.953 | 0.940 |
| **Macro avg** | 0.922 | 0.929 | 0.915 | 0.945 | 0.950 | 0.941 |
| **Detection** | 0.932 | 0.941 | 0.924 | 0.962 | 0.969 | 0.956 |

Table 2: Model performance metrics on the test set. Prec = precision; Avg = average; Detection = Matching the text span regardless of the assigned label.

Figure 1 shows the confusion matrix for overlapping boundary matching. 3.0% of model detections were spurious, while 4.4% of ground truth spans were not detected. Of all spurious classifications, 36.8% were symptoms. The model missed 6.2% of symptoms and 5.6% of results.

### 3.2 Clinical Utility Evaluation

The model identified the diagnosis and at least one relevant symptom in all 257 patients (100% detection). Relevant treatments were detected in 99.2% of patients, missing only two cases: epilepsy ("at se an"—wait and see) and hernia ("reponere"—reposition/reduction). Diagnostic procedures were identified in 99.6% of cases, with one autism case missing "ADOS" and "WISC" assessment tools.



Figure 1: Confusion matrix for the model on the test set with overlapping boundary matching. O = Non-entity spans.

### 3.3 Bias Evaluation

Appendix Table A1 shows mean, median and range of entity counts by medical area. Psychiatric notes mention more symptoms on average (36) than medical (23) and surgical notes (20), with a wider range (0–237 vs. 0–90 and 0–86, respectively). They include fewer anatomies (7 vs. 11 and 12, respectively) and results (6 vs. 14 and 11, respectively).

Figure 2 shows the bootstrapped 95% CIs for macro and micro averaged F1 scores for comparison inside groups. The CI for children is non-overlapping and lower than for seniors but overlap with adults. The psychiatry CI is non-overlapping and lower than the medical and surgical CIs. The observed differences in means for the non-overlaps are at or below 0.017. All other CIs overlap.

Figure 3 shows the bootstrapped 95% CIs for F1 scores for each entity by group. The CIs for diagnostic, anatomy, and result entities overlap inside all groups. In contrast, for disease and symptom entities, the CIs for psychiatry are non-overlapping and lower than those for medical and surgical. For treatment entities, the psychiatry CI is lower than medical but overlap with surgical. The observed differences in means for the non-overlaps are at or below 0.033.

Mean F1 scores for all group levels range from 0.958 to 1.000, with full details on CIs reported in Appendix Table A2 and A3.

Figure 2: Bootstrapped 95% confidence intervals (CIs) for macro (left) and micro (right) averaged F1 scores by group. Note that comparison is only possible between the levels of each group, not between groups.

## 4 Discussion

We saw a substantial improvement in model performance, with exact boundary macro F1 increasing from 0.601 in the original work to 0.922, and overlap boundary F1 from 0.682 to 0.945. Likely causes include more consistent annotation by a single MD rather than six in the original work (Laursen et al., 2023a), a tripled dataset size, improved postprocessing with span merging and voting, and an updated annotation scheme that includes both symptoms, and normal and pathological findings under the Symptom category—reflecting their often similar context in clinical text and simplifying the classification task.

The excellent performance of the prospective utility evaluation of the model shows how it can be used to retrieve all relevant information for physicians managing patients of all ages and genders for the diagnoses included in the evaluation. Given the heterogeneous clinical presentations of the 15 diagnoses evaluated, these results suggest promising potential for broader implementation across all ICD-10 diagnoses.

The findings from our bias study indicate a small but consistent reduction in model performance on psychiatric notes, with minimal effects observed in other groups. While statistical significance was not formally assessed, these differences likely stem from the distinct structure and content of psychiatric notes. Based on clinical experience, psychiatric notes tend to be longer. They also mention more symptoms and contain fewer references to anatomy, results, and diagnostic tests (see Appendix Table A1). These factors suggest that the model could benefit from additional psychiatric notes in training data, although the current performance differences remain very small.

## 5 Conclusion

We present Dora, a Danish clinical NER model that identifies key clinical entities: diseases, symptoms (including normal and pathological findings), diagnostics, treatments, anatomies, and results. Dora achieves substantial improvements over the original model, with a macro F1 score of 0.922 for exact boundary matching and 0.945 for overlapping boundaries. Prospective utility evaluation demonstrates excellent performance in extracting relevant information for physicians. Our bias study reveals a small but consistent performance reduction on psychiatric notes, with minimal variation in other groups, though overall differences remain very small.

## Limitations

While the bias analysis offers valuable insights, several limitations remain; firstly, using four notes per patient may not fully represent the medical condition in case of complex or chronic illnesses. To address this, we applied conservative metrics and bootstrapping in order to improve robustness. Secondly, ground truth labels for the bias study were

Figure 3: Bootstrapped 95% confidence intervals (CIs) for F1 scores for each entity by group. Note that comparison is only possible between the levels of each group, not between groups.

created by correcting model predictions, which, while efficient, may have influenced the annotations. Finally, notes were randomly sampled from a 1.5-year period during which the patient was given the relevant diagnosis. This approach ensures a diverse range of note types per medical area, improving generalisability. It may, however, introduce noise as some notes risk not being strongly representative of their originating medical area.

We cannot rule out that some individual sentences from the evaluation cohort may also appear in the training data. However, since evaluation was performed on full EHRs, the presence of single duplicate sentences, which are common due to standard phrasing in EHRs, is unlikely to impact results.

## Ethics Statement

This study was conducted using clinical data accessed with appropriate institutional permissions. All data usage complied with relevant ethical guidelines and data protection regulations, and was approved by the data providers.

The dataset and model are not publicly available due to sensitive content. Please contact us for sharing options.

## References

AI Sweden. RoBERTa Large. Accessed: 2025-07-30.

Emily Alsentzer, John R Murphy, Willie Boag, Wei-Hung Weng, Di Jin, Tristan Naumann, WA Redmond, and Matthew BA McDermott. 2019. Publicly available clinical bert embeddings. *NAACL HLT 2019*, page 72.

Nancy Chinchor and Beth M Sundheim. 1993. Muc-5 evaluation metrics. In *Fifth Message Understanding Conference (MUC-5): Proceedings of a Conference Held in Baltimore, Maryland, August 25-27, 1993*.

Kevin Clark, Minh-Thang Luong, Quoc V. Le, and Christopher D. Manning. 2020. ELECTRA: pre-training text encoders as discriminators rather than generators. In *8th International Conference on Learning Representations, ICLR 2020, Addis Ababa, Ethiopia, April 26-30, 2020*. OpenReview.net.

Danish Health Data Authority. n.d. The danish medical coding classification system. Accessed: 2025-06-05.

Marzyeh Ghassemi, Tristan Naumann, Peter Schulam, Andrew L Beam, Irene Y Chen, and Rajesh Ranganath. 2020. A review of challenges and opportunities in machine learning for health. *AMIA Summits on Translational Science Proceedings*, 2020:191.

Min Jiang, Yukun Chen, Mei Liu, S Trent Rosenbloom, Subramani Mani, Joshua C Denny, and Hua Xu. 2011. A study of machine-learning-based approaches to extract clinical entities and their assertions from discharge summaries. *Journal of the American Medical Informatics Association*, 18(5):601–606.

Christopher J Kelly, Alan Karthikesalingam, Mustafa Suleyman, Greg Corrado, and Dominic King. 2019. Key challenges for delivering clinical impact with artificial intelligence. *BMC medicine*, 17:1–9.

Martin Laursen, Jannik Pedersen, Rasmus Hansen, Thiusius Rajeeth Savarimuthu, and Pernille Vinholt. 2023a. Danish clinical named entity recognition and relation extraction. In *Proceedings of the 24th Nordic Conference on Computational Linguistics (NoDaLiDa)*, pages 655–666.

Martin S Laursen, Jannik S Pedersen, Rasmus S Hansen, Thiusius R Savarimuthu, Rasmus B Lynggaard, and Pernille J Vinholt. 2023b. Doctors identify hemorrhage better during chart review when assisted by artificial intelligence. *Applied clinical informatics*, 14(04):743–751.

Jinhyuk Lee, Wonjin Yoon, Sungdong Kim, Donghyeon Kim, Sunkyu Kim, Chan Ho So, and Jaewoo Kang. 2020. Biobert: a pre-trained biomedical language representation model for biomedical text mining. *Bioinformatics*, 36(4):1234–1240.

Ilya Loshchilov and Frank Hutter. 2019. Decoupled weight decay regularization. In *International Conference on Learning Representations*.

Jannik S Pedersen, Martin S Laursen, Cristina Soguero-Ruiz, Thiusius R Savarimuthu, Rasmus Søgaard Hansen, and Pernille J Vinholt. 2022. Domain over size: Clinical electra surpasses general bert for bleeding site classification in the free text of electronic health records. In *2022 IEEE-EMBS International Conference on Biomedical and Health Informatics (BHI)*, pages 1–4. IEEE.

Peng Qi, Yuhao Zhang, Yuhui Zhang, Jason Bolton, and Christopher D. Manning. 2020. Stanza: A Python natural language processing toolkit for many human languages. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics: System Demonstrations*.

Pranav Rajpurkar, Emma Chen, Oishi Banerjee, and Eric J Topol. 2022. Ai in health and medicine. *Nature medicine*, 28(1):31–38.

Maria Skeppstedt, Maria Kvist, Gunnar H Nilsson, and Hercules Dalianis. 2014. Automatic recognition of disorders, findings, pharmaceuticals and body structures from clinical text: An annotation and machine learning study. *Journal of biomedical informatics*, 49:148–158.

Ewout W Steyerberg, Frank E Harrell Jr, Gerard JJM Borsboom, MJC Eijkemans, Yvonne Vergouwe, and

J Dik F Habbema. 2001. Internal validation of predictive models: efficiency of some procedures for logistic regression analysis. *Journal of clinical epidemiology*, 54(8):774–781.

Özlem Uzuner, Brett R South, Shuying Shen, and Scott L DuVall. 2011. 2010 i2b2/va challenge on concepts, assertions, and relations in clinical text. *Journal of the American Medical Informatics Association*, 18(5):552–556.

Thomas Vakili, Martin Hansson, and Aron Henriksson. 2025. Sweclineval: A benchmark for swedish clinical natural language processing. In *Proceedings of the Joint 25th Nordic Conference on Computational Linguistics and 11th Baltic Conference on Human Language Technologies (NoDaLiDa/Baltic-HLT 2025)*, pages 767–775.

World Health Organization. 2016. *International Statistical Classification of Diseases and Related Health Problems, 10th Revision (ICD-10)*. https://icd.who.int/browse10/2016/en.

Yuhao Zhang, Yuhui Zhang, Peng Qi, Christopher D Manning, and Curtis P Langlotz. 2021. Biomedical and clinical english model packages for the stanza python nlp library. *Journal of the American Medical Informatics Association*, 28(9):1892–1899.

Zexuan Zhong and Danqi Chen. 2021. A frustratingly easy approach for entity and relation extraction. In *Proceedings of the 2021 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pages 50–61.

# A Bias Evaluation: Detailed Results

|  |  | Mean | Median | Range |
|---|---|---|---|---|
| **Disease** | | | | |
| | medical | 4.66 | 3 | (0 - 27) |
| | psychiatry | 5.82 | 3 | (0 - 41) |
| | surgical | 3.57 | 3 | (0 - 18) |
| **Symptom** | | | | |
| | medical | 22.57 | 18 | (0 - 90) |
| | psychiatry | 35.98 | 19 | (0 - 237) |
| | surgical | 19.77 | 17 | (0 - 86) |
| **Treatment** | | | | |
| | medical | 12.96 | 10 | (0 - 45) |
| | psychiatry | 13.04 | 10 | (0 - 87) |
| | surgical | 15.24 | 12 | (0 - 50) |
| **Diagnostic** | | | | |
| | medical | 21.16 | 19 | (1 - 78) |
| | psychiatry | 12.10 | 8 | (0 - 83) |
| | surgical | 16.52 | 15 | (0 - 50) |
| **Anatomy** | | | | |
| | medical | 11.02 | 7.5 | (0 - 46) |
| | psychiatry | 6.96 | 3 | (0 - 47) |
| | surgical | 12.06 | 9 | (0 - 71) |
| **Result** | | | | |
| | medical | 13.72 | 12 | (0 - 58) |
| | psychiatry | 6.02 | 3 | (0 - 31) |
| | surgical | 11.04 | 8 | (0 - 39) |

Table A1: Mean, median and range of entity counts by medical area for the bias evaluation.

|  | Macro F1 | Micro F1 |
| --- | --- | --- |
| **Gender** | | |
| female | 0.991 (0.986 - 0.995) | 0.991 (0.987 - 0.994) |
| male | 0.989 (0.982 - 0.995) | 0.994 (0.992 - 0.996) |
| **Age group** | | |
| child | 0.985 (0.977 - 0.993) | 0.988 (0.983 - 0.993) |
| adult | 0.989 (0.981 - 0.996) | 0.992 (0.989 - 0.995) |
| senior | 0.996 (0.994 - 0.997) | 0.996 (0.994 - 0.997) |
| **Medical area** | | |
| medical | 0.995 (0.993 - 0.997) | 0.995 (0.994 - 0.997) |
| psychiatry | 0.979 (0.968 - 0.989) | 0.985 (0.980 - 0.990) |
| surgical | 0.996 (0.994 - 0.998) | 0.996 (0.994 - 0.998) |

Table A2: Bootstrapped macro and micro F1 scores with 95% confidence intervals reported by group levels for the bias evaluation.

|  | Disease | Symptom | Treatment |
| --- | --- | --- | --- |
| **Gender** | | | |
| female | 0.994 (0.990 - 0.998) | 0.989 (0.984 - 0.993) | 0.996 (0.993 - 0.998) |
| male | 0.981 (0.964 - 0.998) | 0.993 (0.990 - 0.995) | 0.995 (0.991 - 0.998) |
| **Age group** | | | |
| child | 0.969 (0.942 - 0.996) | 0.986 (0.980 - 0.992) | 0.993 (0.988 - 0.998) |
| adult | 0.997 (0.993 - 1.000) | 0.991 (0.987 - 0.995) | 0.996 (0.992 - 0.999) |
| senior | 0.996 (0.993 - 0.999) | 0.994 (0.991 - 0.997) | 0.996 (0.993 - 0.999) |
| **Medical area** | | | |
| medical | 0.999 (0.997 - 1.000) | 0.994 (0.990 - 0.998) | 1.000 (0.999 - 1.000) |
| psychiatry | 0.966 (0.938 - 0.993) | 0.982 (0.976 - 0.988) | 0.989 (0.983 - 0.995) |
| surgical | 0.997 (0.994 - 1.000) | 0.995 (0.992 - 0.997) | 0.996 (0.993 - 0.999) |

|  | Diagnostic | Anatomy | Result |
| --- | --- | --- | --- |
| **Gender** | | | |
| female | 0.992 (0.985 - 0.998) | 0.984 (0.969 - 0.999) | 0.987 (0.978 - 0.996) |
| male | 0.995 (0.992 - 0.998) | 0.982 (0.965 - 0.999) | 0.977 (0.960 - 0.994) |
| **Age group** | | | |
| child | 0.988 (0.978 - 0.997) | 0.995 (0.991 - 0.999) | 0.967 (0.939 - 0.994) |
| adult | 0.996 (0.992 - 0.999) | 0.964 (0.929 - 0.998) | 0.991 (0.985 - 0.997) |
| senior | 0.997 (0.994 - 0.999) | 0.998 (0.996 - 1.000) | 0.989 (0.979 - 0.999) |
| **Medical area** | | | |
| medical | 0.996 (0.992 - 0.999) | 0.993 (0.985 - 1.000) | 0.985 (0.975 - 0.994) |
| psychiatry | 0.989 (0.980 - 0.997) | 0.958 (0.917 - 0.998) | 0.966 (0.935 - 0.997) |
| surgical | 0.995 (0.989 - 1.000) | 0.997 (0.994 - 1.000) | 0.994 (0.988 - 0.999) |

Table A3: Bootstrapped F1 scores and 95% confidence intervals by entity type and group levels.

## B   Clinical Utility Evaluation: Expected Findings

| MEDICAL | | | |
|---|---|---|---|
| **Disease** | **Symptom** | **Diagnostic** | **Treatment** |
| Epilepsia [DG40] | Seizures<br>Impaired consciousness<br>Tongue bite<br>Urination<br>Amnesia | Blood samples<br>Imaging<br>Electro-<br>encephalogram | Antiseizure medicine |
| Asthma [DJ45] | Dyspnoea<br>Cough | Pulse Oximetry<br>Imaging<br>Blood samples<br>a-puncture<br>pH<br>Pulmonary function test | Bronchodilator<br>Oxygen<br>Steroid |
| Diabetic ketoacidosis type 1 [DE101] | Polyuria/polydipsia<br>Respiratory changes<br>Nausea/vomiting<br>Foetor ex ore<br>Abdominal pain<br>Weakness/fatigue<br>Impaired consciousness | Blood samples<br>a-puncture<br>Glucose<br>Urine sample | Insulin<br>Fluid therapy |
| Rheumatoid arthritis [DM05, DM08] | Pain<br>Swelling<br>Redness<br>Heat of joint(s)<br>Fever<br>Fatigue<br>Other systemic symptoms | Blood samples<br>Imaging | Anti-inflammatory drugs<br>Immunomodulatory drugs<br>Analgesics |
| Pneumonia [DJ189] | Dyspnoea<br>Cough<br>Fever | Pulse Oximetry<br>Imaging<br>Blood samples<br>a-puncture<br>pH | Bronchodilator<br>Oxygen<br>Steroid<br>Antibiotics<br>Fluid therapy |

Table B1: Expected clinical findings in the health record for each medical diagnosis by entity type.

| PSYCHIATRIC | | | |
|---|---|---|---|
| **Disease** | **Symptom** | **Diagnostic** | **Treatment** |
| Generalized anxiety [DF411] | Anxiety<br>Headache<br>Restlessness<br>Pain<br>Tension<br>Fear<br>Sleep disturbances<br>Autonomic hyperactivity<br>Tension | Psychiatric assessment | Psychotherapy<br>Antidepressants<br>CNS depressants |
| Depression [DF33] | Poor concentration<br>Feelings of excessive<br>guilt or low self-worth<br>Hopelessness<br>Thoughts about dying<br>or suicide<br>Disrupted sleep<br>Changes in appetite<br>or weight<br>Feeling very tired<br>or low in energy | MDI or Hamilton scale<br>Psychiatric assessment | Antidepressants<br>Psychotherapy<br>Sleep medication |
| Autism [DF840] | Deficits within:<br>Communication,<br>interaction,<br>and behaviour | Psychiatric assessment | Psychotherapy<br>Sleep medication |
| Suicide attempt /self-injury [DZ915A] | Intentional cause of<br>injury on oneself | Psychiatric assessment | Psychotherapy<br>Antipsychotics<br>CNS depressants |
| Eating disorder [DF50] | Disturbance in one's<br>eating behaviors that<br>affect the person's<br>physical or mental health | Psychiatric assessment<br>BMI | Psychotherapy<br>Enteral/parenteral<br>nutrition therapy |

Table B2: Expected clinical findings in the health record for each psychiatric diagnosis by entity type.

| SURGICAL | | | |
|---|---|---|---|
| **Disease** | **Symptom** | **Diagnostic** | **Treatment** |
| Appendicitis [DK35, DK37, DK379] | Pain Fever Nausea/vomiting | Abdominal examination Blood samples Imaging | Surgery/appendectomy Antibiotics Analgesics |
| Epistaxis [DR040C, DR040A, DR040B] | Bleeding from nose nose or mouth | Blood samples Rhino endoscopy Imaging | Compressive therapy Ablation Haemostatics Transfusion Fluid therapy |
| Fracture of tibia [DS821, DS823] | Pain Swelling Loss of function Displacement | Examination Imaging | Analgesics Fixation Surgery Antibiotics |
| Hernia [DK409] | Pain Protrusion Fever Nausea/vomiting | Abdominal examination Imaging Blood samples | Surgery Antibiotics Analgesics |
| Ilieus [DK567] | Pain Fever Nausea/vomiting | Abdominal examination Imaging Blood samples | Surgery Antibiotics Analgesics |

Table B3: Expected clinical findings in the health record for each surgical diagnosis by entity type.

# CMC-SC: Cross-Modal Contextualized ASR Spelling Correction via BERT and WavLM using a Soft Fusion Framework

**Mohammad Reza Peyghan** and **Sajjad Amini**[*] and **Shahrokh Ghaemmaghami**
Electronics Research Institute and Department of Electrical Engineering
Sharif University of Technology
{m.peyghan, s_amini, ghaemmag}@sharif.edu

## Abstract

Automatic Speech Recognition (ASR) systems remain error-prone in challenging acoustic conditions, leading to spelling mistakes that degrade downstream applications. Despite the surge in the number of studies on post-refinement methods, existing Spelling Correction (SC) approaches often rely solely on textual cues or phonetic features, limiting their ability to provide speech-aware corrections. In this work, we introduce a Cross-Modal Contextualized Spelling Correction framework (CMC-SC) that jointly incorporates contextualized acoustic and textual information. Unlike prior methods that use phonetics solely for candidate selection, our solution leverages contextualized speech tokens in the generation of corrections, improving accuracy and context awareness. CMC-SC features a detection module to identify errors, a cross-modal correction module to generate fixes using acoustic and textual tokens, and a soft fusion step to refine corrections while retaining context. The proposed method improves error rates compared to baselines and, with only 140M trainable parameters, offers an efficient solution for ASR spelling correction.

## 1 Introduction

Automatic Speech Recognition (ASR) systems have become increasingly important in recent years, enabling a wide range of applications, from virtual assistants to transcription services. The field has seen significant growth, driven by advancements in deep learning and natural language processing. However, despite these advances, ASR systems still face challenges, particularly in diverse acoustic environments and with speakers of different accents Errattahi et al. (2018). Retraining ASR

models with domain-specific data can often mitigate these issues to some extent, but in many cases, the ASR model is not accessible for direct modification, functioning as a black box. In such scenarios, post-refinement techniques can be effectively employed to improve transcription quality.

Various ASR refinement techniques have been explored, especially since the advent of Transformers Vaswani et al. (2017). Broadly speaking, ASR refinement can be categorized into three main classes: fusion, re-scoring, and correction.

Fusion methods aim to improve ASR first-pass decoding by integrating external linguistic information at each decoding step. These techniques typically augment the ASR decoder's internal language model with external Language Models (LMs), whether a simple n-gram Kannan et al. (2018), a neural LM Kim et al. (2021), or a Large Language Model (LLM) Hori et al. (2025).

The Re-Scoring paradigm, by contrast, is a second-pass scheme that assumes the 1-best ASR hypothesis may not properly represent the information from the decoding step. This paradigm generates an $N$-best list of hypotheses and uses an external model (e.g., an n-gram or neural language model) to re-rank those candidates, selecting a linguistically superior candidate Shin et al. (2019); Gandhe and Rastrow (2020).

Correction approaches tackle the problem by revising a given ASR transcript to produce a new, improved sequence. Some correction techniques employ a second-pass decoding strategy, where a second decoder (or encoder-decoder) reconsiders acoustic features or the initial hypothesis. This decoding step can utilize an n-gram Bassil and Semaan (2012), a neural LM Zhang et al. (2019), or an LLM Udagawa et al. (2024), whether adopting both modalities Orihashi et al. (2021); Xia et al. (2017) or text-only correction Hrinchuk et al. (2020); Jia et al. (2025). In recent research, researchers have used Retrieval Augmented Gen-

---

Figure 1: The Overall Diagram of the Proposed CMC-SC.

eration (RAG) with an external corpus for transcript correction Robatian et al. (2025); Gong et al. (2025).

Another notable approach to ASR error correction is the use of encoder transformers, primarily BERT Devlin et al. (2019). These models leverage their contextual understanding to replace erroneous tokens. However, using pre-trained BERT alone is suboptimal for ASR correction Zhang et al. (2020) due to (1) reliance on textual cues, which risks incorrect substitutions, and (2) a domain mismatch between its clean pre-training data and noisy ASR outputs. To address this, FASPell Hong et al. (2019) employs a Confidence-Similarity Decoder (CSD) to filter BERT's candidates by phonetic and orthographic similarity. Similarly, SpellGCN Cheng et al. (2020) enhances BERT with a Graph Convolutional Network (GCN) to model phonological and symbolic relations. Other methods incorporate detection modules. For instance, a method Zhang et al. (2020) detects and softly masks probable errors based on confidence scores, feeding them into a correction model and summing outputs with original embeddings. Another approach Zhang et al. (2021) fuses token and phonetic embeddings post-detection for phonetic-aware correction. Additionally, a dynamic error scaling method Fan et al. (2023) integrates words and pinyin for semantically and phonetically aware character-level correction.

However, encoder-based methods addressing these challenges often rely on phonetic information derived from text, which can be misleading. Our method addresses this issue by:

1. Extracting contextualized acoustic information directly from speech using WavLM Chen et al. (2022), unlike Fan et al. (2023) and Zhang et al. (2021), which rely on transcription-based information.

2. Joint processing of contextualized acoustic and textual tokens through a Cross-Modal BERT (CM-BERT) unlike Hong et al. (2019); Cheng et al. (2020), which rely on phonetic information in a secondary branch.

3. Using a soft fusion technique to combine CM-BERT outputs with original token embeddings, preserving transcription information, unlike the direct summing approach used in Zhang et al. (2021).

Finally, our approach improves upon existing baselines by a large margin, demonstrating its effectiveness in improving ASR quality.

## 2 Method

This section presents the methodology for enhancing ASR transcriptions using a cross-modal framework that integrates textual and acoustic data. The approach comprises two main components: a detection module to identify erroneous tokens and a cross-modal correction module to rectify these errors using a soft-fusion framework. The structure of the proposed CMC-SC is illustrated in Figure 1, and subsequent subsections detail each component.

### 2.1 Data Pre-Processing

We perform the following data preprocessing steps to enable end-to-end (E2E) training of our proposed model:

1. We run paired speech-text examples through a black-box ASR to obtain its transcriptions.

2. For each utterance, we align the ASR transcription with the corresponding ground-truth using Levenshtein alignment (edit distance).

272

From this, we create a per-token binary sequence where '1' indicates an erroneous token and '0' indicates a correctly recognized token.

Consequently, our dataset comprises samples in the form of (utterance, erroneous text, labels, target text). The (erroneous text, labels) pairs are used to train the Detection module, while the (utterance, erroneous text, target text) pairs are used to train the Correction module.

## 2.2 Detection and Masking

The detection module, the first stage of our spelling correction pipeline, aims to identify erroneous tokens in ASR transcriptions to enable targeted corrections and prevent over-correction, which is a common issue in E2E methods Imai et al. (2025). It integrates a frozen and pre-trained BERT model (denoted T-BERT) to extract contextual token embeddings, combines these embeddings with initial token embeddings via a residual connection, and feeds them to a two-layer BiLSTM classification head. The residual connection is crucial, as it allows the model to consider content other than context, addressing potential misclassifications from incorrect tokens affecting the frozen BERT's embeddings. Then, a linear layer outputs logits for each token, indicating whether it is erroneous. This module is trained using Binary Cross-Entropy (BCE) loss with logits and per-token binary labels:

$$\text{BCE}_D = \frac{1}{T} \sum_{i=1}^{T} \log\big(1 + e^{-(2y_i - 1)z_i}\big).$$

where $z_i$ is the logit for token $i$, and $y_i \in \{0, 1\}$ indicates if the token is erroneous.

At inference time, we compute the sigmoid of each logit and compare it to a predefined threshold; tokens with likelihood above this threshold are deemed erroneous and replaced with the [MASK] token. This masking strategy ensures that CM-BERT's context is derived from the most probable correct tokens, preventing incorrect tokens from negatively affecting the contextual representation.

## 2.3 Soft Fusion and Cross-Modal Correction

The cross-modal correction module refines ASR transcriptions by integrating textual and acoustic data to produce accurate, speech-aware, and contextually appropriate corrections. Using both modalities, it improves transcription quality using a cross-modal and joint attention approach.

The correction module receives a sequence of token embeddings from the detection phase, where tokens identified as incorrect are replaced with the [MASK] token, denoted as $\mathbf{E}_m$. It also extracts contextualized speech features from raw audio using a pre-trained WavLM network. These speech features are then projected to match the dimensionality of the CM-BERT, resulting in $\mathbf{E}_S$.

The masked text embeddings $\mathbf{E}_m$ and the projected speech embeddings $\mathbf{E}_S$ are concatenated to form the input $\mathbf{E}_{\text{in}} = [\mathbf{E}_m; \mathbf{E}_S]$. This concatenated input is then processed by the CM-BERT, a transformer-based model that outputs contextualized, speech-aware representations by enabling cross-modal interactions between text and speech through its attention mechanisms.

To prevent over-correction and preserve correct tokens, the Soft-Fusion (SF) strategy blends each token's original embedding $\mathbf{E}^{(i)}$ with its corresponding cross-modal contextual embedding $\mathbf{E}_c^{(i)}$, based on a confidence score $\alpha_i$ from the detection phase that indicates the likelihood that token $i$ is incorrect. Specifically, under the SF strategy, the output embedding for each token $i$ is computed as:

$$\mathbf{E}_o^{(i)} = (1 - \alpha_i) \cdot \mathbf{E}^{(i)} + \alpha_i \cdot \mathbf{E}_c^{(i)}$$

As a result, tokens with a low $\alpha_i$ (indicating they are likely correct) retain more of their original embedding, while tokens with a high $\alpha_i$ (indicating they are likely incorrect) incorporate more of the speech-informed representation. This adaptive interpolation ensures precise corrections where needed while preserving accurate text.

Finally, the softly-fused embeddings are classified into tokens using a softmax layer, guided by the Cross Entropy (CE) loss with logits and token IDs. The CE loss is given by:

$$\text{CE}_C = -\frac{1}{T} \sum_{t=1}^{T} \log\left(\frac{\exp(z_{t,y_t})}{\sum_{k=1}^{V} \exp(z_{t,k})}\right)$$

where $T$ is the sequence length, $V$ is the vocabulary size, $z_{t,k}$ is the logit for token $k$ at position $t$, and $y_t$ is the true token ID.

## 3 Experiments and Results

In this section, we detail the experiments conducted to assess the proposed model and present the results in comparison to several baseline models, which were re-implemented to ensure a fair evaluation. Additionally, we assess the performance of each

| | Model | Parameters | | Detection | | | Correction | | | Error Rate | |
|---|---|---|---|---|---|---|---|---|---|---|---|
| | | Total | Trainable | P | R | F1 | P | R | F1 | Word | Character |
| **Comparative Study** | Whisper-Tiny (Baseline) Radford et al. (2023) | 39M | - | - | - | - | - | - | - | 24.5 | 17.2 |
| | Whisper-Small Radford et al. (2023) | 244M | - | - | - | - | - | - | - | 13.7 | 6.1 |
| | Whisper-Medium Radford et al. (2023) | 769M | - | - | - | - | - | - | - | 11.7 | **4.2** |
| | PT-BERT+BiLSTM (Multi-Task Training) | 140M | 30M | 85.96 | 85.90 | 85.88 | 81.74 | 83.94 | 81.82 | 18.2 | 17.5 |
| | FT-BERT+BiLSTM (Multi-Task Training) | 140M | 140M | 85.84 | 85.95 | 85.79 | 82.08 | 85.77 | 82.18 | 17.9 | 17.1 |
| | Soft-Masked BERT Zhang et al. (2020) | 250M | 250M | 86.14 | 86.23 | 86.12 | 87.85 | 87.23 | 87.18 | 13.2 | 9.8 |
| | CMC (Ours) | 300M | 140M | **87.32** | **87.52** | **87.30** | **91.31** | **91.35** | **91.27** | **9.2** | 5.6 |
| **Ablation Study** | CMC – WavLM | 210M | 140M | 87.18 | 87.24 | 87.15 | 88.91 | 88.74 | 88.64 | 12.1 | 9.6 |
| | CMC – SF | 300M | 140M | 87.17 | 87.18 | 87.15 | 89.03 | 89.12 | 89.01 | 10.7 | 7.5 |

Table 1: Comparative and Ablation Studies (all refinement methods are applied to Whisper-Tiny)

| Transcriptions | 1) the cut start on the fence<br>2) she begged home smiling all the way knowing that she had won |
|---|---|
| Detection Predictions | 1) [1, 1, 1, 0, 0, 0]<br>2) [0, 1, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0] |
| Detection Labels | 1) [0, 1, 1, 0, 0, 0]<br>2) [0, 1, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0] |
| Refined | 1) the cat sat on the fence<br>2) she biked home smiling all the way knowing that she had won |
| Ground Truths | 1) the cat sat on the fence<br>2) she biked home smiling all the way knowing that she had won |

Table 2: Examples of CMC-SC on the Common Voice test set.

module within the model through an ablation study, systematically removing each module to evaluate its impact on the overall performance.

To evaluate our model, we introduce baseline models. We use three ASR models (Whisper-Tiny, Whisper-Small, Whisper-Medium) to assess the importance of post-refinement and Cross-Modal attention against adopting larger ASR systems. We also trained two spelling correction baselines, Pre-Trained (PT) and Fine-Tuned (FT) BERT, following Zhang et al. (2020); Cheng et al. (2020); Fan et al. (2023), to highlight our model's contribution. Plus, we re-implemented Soft-Masked BERT Zhang et al. (2020) as another benchmark.

We perform an ablation study to quantify each module's contribution to the CMC-SC model. First, we remove the speech tokens (i.e., contextualized acoustic information) and retrain under identical conditions, noting that CM-BERT is originally pre-trained on text, so its performance may still reflect textual bias rather than a true absence of cross-modal data. This ablation also underscores the significance of the residual connection in the detection module, which is the primary distinction of this module in the ablation and compared to the PT-BERT. Next, we remove the Soft-Fusion module, which retains information from the original transcription, and train it again. Table 1 presents these results, demonstrating that each module positively impacts the overall performance of CMC-SC.

All experiments ran on an NVIDIA RTX 3090 GPU for 30 epochs using the AdamW optimizer. The best model uses a batch size of 32, a learning rate of $1 \times 10^{-5}$ with a linear scheduler, both T-BERT and CM-BERT have a maximum context length of 128 tokens, and the threshold in the masking module is set empirically to $0.5$. To align speech tokens with BERT embeddings, we project them into 50 tokens of dimension 768. We have employed the Mozilla Common Voice dataset Ardila et al. (2019) (original train/dev/test splits) and report results on its test set.

Finally, as shown in Table 1, our proposed method improves the baselines by a large margin, demonstrating substantial potential to improve the spelling correction task. Notably, our model has only 140M trainable parameters and outperforms the pre-trained Whisper-medium with 769M parameters, making it a lightweight yet effective solution. The examples of CMC-SC are provided in Table 2.

## 4 Conclusion

In this paper, we have introduced Cross-Modal Contextualized Spelling Correction (CMC-SC), a novel framework designed to enhance ASR transcription accuracy by correcting spelling errors. CMC-SC integrates a detection module using a frozen BERT model and BiLSTM to identify errors by capturing contextual and sequential patterns, and a correction module that blends text embeddings with acoustic features from a pretrained

WavLM. This approach ensures precise, context-aware corrections while preserving accurate tokens via a soft fusion framework. Experiments show CMC-SC reduces error rates with only 140 million trainable parameters, balancing performance and computational efficiency. Future work includes supporting additional languages and integrating advanced pretrained cross-modal networks for deeper linguistic and acoustic insights.

## Limitations

Despite the resulting advancements, ASR models remain error-prone in challenging environments. In clean settings, errors are primarily substitutions or spelling mistakes, for which spelling correction methods are computationally efficient. However, the proposed method may be less effective for errors involving insertions and deletions. Additionally, trained on general data, the model may require re-training for domain-specific applications, such as medical terminology.

## References

Rosana Ardila, Megan Branson, Kelly Davis, Michael Henretty, Michael Kohler, Josh Meyer, Reuben Morais, Lindsay Saunders, Francis M Tyers, and Gregor Weber. 2019. Common voice: A massively-multilingual speech corpus. *arXiv preprint arXiv:1912.06670*.

Youssef Bassil and Paul Semaan. 2012. Asr context-sensitive error correction based on microsoft n-gram dataset. *arXiv preprint arXiv:1203.5262*.

Sanyuan Chen, Chengyi Wang, Zhengyang Chen, Yu Wu, Shujie Liu, Zhuo Chen, Jinyu Li, Naoyuki Kanda, Takuya Yoshioka, Xiong Xiao, et al. 2022. Wavlm: Large-scale self-supervised pre-training for full stack speech processing. *IEEE Journal of Selected Topics in Signal Processing*, 16(6):1505–1518.

Xingyi Cheng, Weidi Xu, Kunlong Chen, Shaohua Jiang, Feng Wang, Taifeng Wang, Wei Chu, and Yuan Qi. 2020. Spellgcn: Incorporating phonological and visual similarities into language models for chinese spelling check. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 871–881.

Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. 2019. Bert: Pre-training of deep bidirectional transformers for language understanding. In *Proceedings of the 2019 conference of the North American chapter of the association for computational linguistics: human language technologies, volume 1 (long and short papers)*, pages 4171–4186.

Rahhal Errattahi, Asmaa El Hannani, and Hassan Ouahmane. 2018. Automatic speech recognition errors detection and correction: A review. *Procedia Computer Science*, 128:32–37.

Jiaxin Fan, Yong Zhang, Hanzhang Li, Jianzong Wang, Zhitao Li, Sheng Ouyang, Ning Cheng, and Jing Xiao. 2023. Boosting chinese asr error correction with dynamic error scaling mechanism. In *Proc. Interspeech 2023*, pages 2173–2177.

Ankur Gandhe and Ariya Rastrow. 2020. Audio-attention discriminative language model for asr rescoring. In *ICASSP 2020-2020 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, pages 7944–7948. IEEE.

Xun Gong, Anqi Lv, Zhiming Wang, Huijia Zhu, and Yanmin Qian. 2025. Br-asr: Efficient and scalable bias retrieval framework for contextual biasing asr in speech llm. *arXiv preprint arXiv:2505.19179*.

Yuzhong Hong, Xianguo Yu, Neng He, Nan Liu, and Junhui Liu. 2019. Faspell: A fast, adaptable, simple, powerful chinese spell checker based on dae-decoder paradigm. In *Proceedings of the 5th Workshop on Noisy User-generated Text (W-NUT 2019)*, pages 160–169.

Takaaki Hori, Martin Kocour, Adnan Haider, Erik McDermott, and Xiaodan Zhuang. 2025. Delayed fusion: Integrating large language models into first-pass decoding in end-to-end speech recognition. *arXiv preprint arXiv:2501.09258*.

Oleksii Hrinchuk, Mariya Popova, and Boris Ginsburg. 2020. Correction of automatic speech recognition with transformer sequence-to-sequence model. In *Icassp 2020-2020 ieee international conference on acoustics, speech and signal processing (icassp)*, pages 7074–7078. IEEE.

Saki Imai, Tahiya Chowdhury, and Amanda Stent. 2025. Evaluating open-source asr systems: Performance across diverse audio conditions and error correction methods. In *Proceedings of the 31st International Conference on Computational Linguistics*, pages 5027–5039.

Linzhao Jia, Han Sun, Yuang Wei, Changyong Qi, and Xiaozhe Yang. 2025. Epic: Error pattern informed correction for classroom asr with limited labeled data. In *ICASSP 2025-2025 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, pages 1–5. IEEE.

Anjuli Kannan, Yonghui Wu, Patrick Nguyen, Tara N Sainath, Zhijeng Chen, and Rohit Prabhavalkar. 2018. An analysis of incorporating an external language model into a sequence-to-sequence model. In *2018 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, pages 1–5828. IEEE.

Suyoun Kim, Yuan Shangguan, Jay Mahadeokar, Antoine Bruguier, Christian Fuegen, Michael L Seltzer, and Duc Le. 2021. Improved neural language model fusion for streaming recurrent neural network transducer. In *ICASSP 2021-2021 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, pages 7333–7337. IEEE.

Shota Orihashi, Ryo Masumura, Mana Ihori, Takafumi Moriya, Akihiko Takashima, Naoki Makishima, Takanori Ashihara, and Tomohiro Tanaka. 2021. Cross-modal transformer-based neural correction models for automatic speech recognition. *Interspeech 2021*.

Alec Radford, Jong Wook Kim, Tao Xu, Greg Brockman, Christine McLeavey, and Ilya Sutskever. 2023. Robust speech recognition via large-scale weak supervision. In *International conference on machine learning*, pages 28492–28518. PMLR.

Amin Robatian, Mohammad Hajipour, Mohammad Reza Peyghan, Fatemeh Rajabi, Sajjad Amini, Shahrokh Ghaemmaghami, and Iman Gholampour. 2025. Gec-rag: Improving generative error correction via retrieval-augmented generation for automatic speech recognition systems. *arXiv preprint arXiv:2501.10734*.

Joonbo Shin, Yoonhyung Lee, and Kyomin Jung. 2019. Effective sentence scoring method using bert for speech recognition. In *Asian Conference on Machine Learning*, pages 1081–1093. PMLR.

Takuma Udagawa, Masayuki Suzuki, Masayasu Muraoka, and Gakuto Kurata. 2024. Robust asr error correction with conservative data filtering. *arXiv preprint arXiv:2407.13300*.

Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. 2017. Attention is all you need. *Advances in neural information processing systems*, 30.

Yingce Xia, Fei Tian, Lijun Wu, Jianxin Lin, Tao Qin, Nenghai Yu, and Tie-Yan Liu. 2017. Deliberation networks: Sequence generation beyond one-pass decoding. *Advances in neural information processing systems*, 30.

Ruiqing Zhang, Chao Pang, Chuanqiang Zhang, Shuohuan Wang, Zhongjun He, Yu Sun, Hua Wu, and Haifeng Wang. 2021. Correcting chinese spelling errors with phonetic pre-training. In *Findings of the Association for Computational Linguistics: ACL-IJCNLP 2021*, pages 2250–2261.

Shaohua Zhang, Haoran Huang, Jicong Liu, and Hang Li. 2020. Spelling error correction with soft-masked bert. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 882–890.

Shiliang Zhang, Ming Lei, and Zhijie Yan. 2019. Investigation of transformer based spelling correction model for ctc-based end-to-end mandarin speech recognition. In *Interspeech*, pages 2180–2184.

# On Limitations of LLM as Annotator for Low Resource Languages

**Suramya Jadhav[1,3], Abhay Shanbhag[1,3], Amogh Thakurdesai[1,3],**
**Ridhima Sinare[1,3], and Raviraj Joshi[2,3]**
[1]Pune Institute of Computer Technology, Pune
[2]Indian Institute of Technology Madras, Chennai
[3]L3Cube Labs, Pune

## Abstract

Low-resource languages face significant challenges due to the lack of sufficient linguistic data, resources, and tools for tasks such as supervised learning, annotation, and classification. This shortage hinders the development of accurate models and datasets, making it difficult to perform critical NLP tasks like sentiment analysis or hate speech detection. To bridge this gap, Large Language Models (LLMs) present an opportunity for potential annotators, capable of generating datasets and resources for these underrepresented languages. In this paper, we focus on Marathi, a low-resource language, and evaluate the performance of both closed-source and open-source LLMs as annotators, while also comparing these results with fine-tuned BERT models. We assess models such as GPT-4o and Gemini 1.0 Pro, Gemma 2 (2B and 9B), and Llama 3.1 (8B and 405B) on classification tasks including sentiment analysis, news classification, and hate speech detection. Our findings reveal that while LLMs excel in annotation tasks for high-resource languages like English, they still fall short when applied to Marathi. Even advanced models like GPT-4o and Llama 3.1 405B underperform compared to fine-tuned BERT-based baselines, with GPT-4o and Llama 3.1 405B trailing fine-tuned BERT by accuracy margins of 10.2% and 14.1%, respectively. This highlights the limitations of LLMs as annotators for low-resource languages.

## 1 Introduction

Even with advancements in NLP, the curation of annotations for supervised tasks like sentiment analysis, text classification, and inference has been the primary responsibility of human linguistic experts (Tan et al., 2024). Data annotations play an integral part in both building and evaluating a model. Hence, the quality and reliability of data lie at the core of the performance and usefulness of the model being built.

The process of curating good-quality data annotations is expensive in terms of time and cost, specifically when it comes to compiling data annotations for low-resource languages. The aim of this study is to explore whether Large Language Models (LLMs) can be effectively leveraged to create supervised data resources for low-resource languages, with Marathi as the focus in this case.

Recent generative models like ChatGPT have shown competitive quality in data annotations for simpler tasks like sentiment analysis while human expert annotations proved to be better for intricate tasks Nasution and Onan (2024). ChatGPT was evaluated by Zhu et al. (2023) to check its capability of reproducing human-generated labels for social computing tasks. In these experiments, ChatGPT obtained an average accuracy of 0.60 with 0.64 being the highest accuracy for the sentiment analysis task. In addition to these, the works of Kuzman et al. (2023); Gao et al. (2023) have previously evaluated ChatGPT's performance with that of human experts. Experiments performed by Mohta et al. (2023) demonstrated that Vicuna 13b performed reasonably well for numerous annotation tasks compared to other models that were tested like Vicuna 7b, Llama (13b, 7b) and Instruct-BLIP(13b, 7b). However, it is important to note that most of these experiments target the English language.

India is a multilingual nation with various regional languages and most of these languages fall under the low-resource (LR) category. Low resource languages are languages such as Marathi and Hindi that lack annotated training datasets and have very few task-specific resources compared to high resource languages such as Spanish and English.

This paper presents a case study on the performance of Large Language Models (LLMs) in annotating the low-resource language Marathi. We conduct a comprehensive comparative analysis of

various closed-source and open-source LLMs, revealing that many LLMs still fall significantly short of the baseline performance achieved by BERT-based models and are not yet capable of replacing human annotators.

Specifically, we evaluated models such as GPT-4o, Gemini 1.0 Pro, Gemma 2 (2B and 9B), Llama 3.1 (8B and 405B) across multiple tasks, including 3-class sentiment analysis, 2-class, and 4-class hate speech detection, as well as news classification based on headlines, long paragraphs, and full documents.

The key contributions of this research work are as follows:

- We have conducted a first-of-its-kind detailed comparative study between fine-tuned BERT models and large language models (LLMs), by evaluating their potential to be used as annotators for a low-resource language, Marathi.

- We observe that the average results of the Few-shot prompting technique outperform the average result of the Zero-shot prompting technique in all the models tested.

- We have provided valuable insights into the effectiveness of both open- and closed-source large language models (LLMs), including GPT-4o, Llama 3.1 405B, Llama 3.1 8B, Gemma 2 9B, Gemma 2 2B, and Gemini 1.0 Pro, on tasks such as Marathi Sentiment Analysis, Hate Speech Detection, and News Categorization. Our results strongly demonstrate that LLMs are still not fully reliable for annotation tasks in Indic languages.

- Model ranking, based on accuracy metrics, is GPT-4o > Llama-3.1-405B > Gemini 1.0 Pro > Gemma 2 9B > Llama 3.1 8B > Gemma-2-2B.

The paper is structured as follows: Section 2 provides a concise review of prior research on data annotation and the use of LLMs. In Section 3, we detail the datasets used and the Section 4 describes models employed in our evaluation. Section 5 describes the experimental setup and the APIs leveraged to evaluate the LLMs. Section 6 presents the results, along with a comparative analysis of various LLMs and BERT-based models, highlighting the key findings of our research. Finally, in Section 7, we conclude our discussion.

## 2 Literature Review

Many low-resource languages, including Marathi, lack well-annotated datasets, making it difficult to train effective models for tasks like sentiment analysis and classification Al-Wesabi et al. (2023). The absence of sufficient data often leads to poor performance in tasks that require labeled corpora R et al. (2023).

Low-resource languages also present unique linguistic challenges not well-represented in high-resource models Krasadakis et al. (2024), highlighting the need for specialized approaches. With the rise of LLMs, these models have been explored as a solution to mitigate the scarcity of annotated data in low-resource languages.

Several works demonstrate the use of LLMs as annotators for low-resource language tasks. Pavlovic and Poesio (2024) reviewed LLMs like GPT-4 and noted performance drops when handling non-English languages. In Kholodna et al. (2024), the authors explored the integration of large language models (LLMs), specifically GPT-4 Turbo, into an active learning framework designed for low-resource language tasks. Their work demonstrates the use of few-shot learning to generate useful annotations, significantly enhancing performance on low-resource tasks. Additionally, they implemented the GPT-4 Turbo model as a classifier within the training loop, leading to a substantial reduction in annotation costs, which were 42.45 times lower compared to traditional methods. However, the general performance of LLMs remains limited, especially for languages with fewer resources Hedderich et al. (2020).

The studies of Ding et al. (2022) and Mohta et al. (2023) further evaluated LLM performance on multilingual datasets, with results indicating that models like GPT-3 and open-source LLMs struggle with non-English data. Srivastava et al. (2022) showed that increasing model size does not consistently enhance performance for low-resource languages, unlike high-resource languages like English.

Bias is another concern with LLMs. Bavaresco et al. (2024) introduced JUDGE-BENCH to evaluate LLM biases, noting that training data heavily influences model outputs, which can be problematic in annotating complex or sensitive tasks in low-resource languages. While LLMs used for high-resource languages (HRL) are giving promising results, that is not the case for low-

resource languages. Nasution and Onan (2024) explored ChatGPT-4's performance in annotation tasks across languages like Turkish and Indonesian, offering insights into LLM applicability for Low Resource Language(LRL), a relevant consideration for our focus on Marathi.

# 3 Dataset

In this research, we focus on three major task categories using relevant Marathi datasets:

1) MahaSent (Kulkarni et al., 2021; Pingle et al., 2023) – classifies sentiment of Marathi tweets into three classes of positive, negative, or neutral categories.

2) MahaHate (Patil et al., 2022) – measures the level of abusive and hostile content in Marathi text. This dataset includes two supervised tasks: MahaHate 2-Class, which categorizes content as either HATE or NOT, and MahaHate 4-Class, which provides finer distinctions with categories: Hate (HATE), Offensive (OFFN), Profane (PRFN), and Not (NOT).

3) MahaNews (Mittal et al., 2023; Mirashi et al., 2024) – classifies headlines and articles from Marathi news. It comprises three supervised datasets: Short Headlines Classification (SHC), Long Document Classification (LDC), and Long Paragraph Classification (LPC), each categorizing news content into 12 classes: Auto, Bhakti, Crime, Education, Fashion, Health, International, Manoranjan, Politics, Sports, Tech, and Travel. The distribution of all the mentioned datasets is provided in Table 2

# 4 Methodolgy

We investigated the distinctions between LLM-generated and human-generated annotations for the Indic language, Marathi, using a comparative methodology, and analyzed the results with fine-tuned BERT-based models for detailed insights.

## 4.1 LLMs

In our annotation experiments, we evaluated the performance of LLMs for the Marathi language using two prompting techniques: zero-shot and few-shot learning. We tested both open-source models (Llama 3.1 8B, Llama 3.1 405B, Gemma 2 2B, and Gemma 2 9B) and closed-source models (Gemini 1.0 Pro, GPT-4o), and compared their results with BERT-based models. The performance

of each LLM under both prompting strategies is summarized in Table 1.

## 4.2 BERT Based Models

We used fine-tuned BERT-based models to compare performance with LLMs, where MahaSent-MD MahaHate-BERT, MahaNews-SHC-BERT, MahaNews-LPC-BERT, and MahaNews-LDC-BERT are fine-tuned versions of MahaBERT, while MahaHate-multi-RoBERTa has MahaRoBERTa as the base model. Each of these models was fine-tuned on the corresponding datasets, and their respective performances are summarized in Table 1.

# 5 Experimental Setup

Our main objective is to assess the LLMs on three different tasks and related datasets to ascertain whether LLMs could take the place of, or at least support, human annotation efforts. We employed both few-shot and zero-shot prompting techniques, with LLM-generated annotations evaluated against the ground truth labels. For all datasets, the test split was used. The open-source models (Llama 3.1 8B, Gemma 2 2B, and Gemma 2 9B) exhibited slower response times and required significant computational resources to generate predictions. However, by utilizing NVIDIA NIM APIs, we were able to accelerate predictions from these models, improving both speed and precision. For the closed-source Gemini 1.0 Pro model, we used the Gemini API, while GPT-4o predictions were generated manually via ChatGPT's default settings to annotate the samples. In our research, we could only use a subset of samples from each dataset due to the restrictive usage regulations and cost limits of the mentioned APIs. To maintain consistency and fairness in the performance comparison, all results from both LLM-based and BERT-based models were evaluated on a uniform subset. Specifically, we evaluated 490 samples from the MahaSent and MahaHate datasets, while for MahaNews, we selected 40 samples from each of the 12 classes, amounting to a total of 480 samples.

# 6 Result

This section provides a detailed overview of the experiments conducted for the annotation of three distinct tasks, utilizing six large language models (LLMs) and six BERT-based models (BERT model fine-tuned on target task). Table 1 summarizes the performance metrics of the fine-tuned BERT-based

| Dataset | Tech | Llama 3.1 8B | Gemma 2 2B | Gemma 2 9B | Gemini 1.0 Pro | Llama 3.1 405B | GPT-4o | Fine Tuned BERT |
|---|---|---|---|---|---|---|---|---|
| MahaSent | ZS | 0.76 | 0.71 | 0.69 | 0.78 | 0.77 | 0.79 | 0.80 |
| | FS | 0.79 | 0.76 | 0.78 | 0.76 | 0.81 | **0.82** | |
| MahaHate-2C | ZS | 0.64 | 0.71 | 0.78 | 0.74 | 0.77 | 0.80 | **0.91** |
| | FS | 0.78 | 0.72 | 0.82 | 0.72 | 0.82 | 0.82 | |
| MahaHate-4C | ZS | 0.40 | 0.39 | 0.43 | 0.43 | 0.49 | 0.58 | **0.73** |
| | FS | 0.48 | 0.41 | 0.46 | 0.45 | 0.52 | 0.60 | |
| MahaNews-SHC | ZS | 0.60 | 0.54 | 0.68 | 0.68 | 0.75 | 0.78 | **0.85** |
| | FS | 0.66 | 0.54 | 0.68 | 0.70 | 0.74 | 0.78 | |
| MahaNews-LPC | ZS | 0.66 | 0.55 | 0.71 | 0.72 | 0.76 | 0.77 | **0.89** |
| | FS | 0.67 | 0.50 | 0.72 | 0.74 | 0.76 | 0.75 | |
| MahaNews-LDC | ZS | 0.69 | 0.62 | 0.78 | 0.74 | 0.76 | 0.81 | **0.96** |
| | FS | 0.69 | 0.62 | 0.80 | 0.75 | 0.78 | 0.81 | |
| Average | ZS | 0.625 | 0.587 | 0.678 | 0.682 | 0.716 | 0.755 | **0.857** |
| | FS | 0.678 | 0.592 | 0.710 | 0.687 | 0.738 | 0.763 | |

Table 1: Model Comparison across different tasks. Tech: Different Prompting Techniques Used; ZS: Zero Shot; FS: Few Shot; 2C: 2-Class; 4C: 4-Class; SHC: Short Headlines Classification; LDC: Long Document Classification; LPC: Long Paragraph Classification; BERT: Refer Section 4.2 for details about BERT models.

| Split | MahaSent | MahaHate 2-C | MahaHate 4-C | SHC | LDC | LPC |
|---|---|---|---|---|---|---|
| Train | 12114 | 30000 | 21500 | 22014 | 22014 | 42870 |
| Valid | 1500 | 3750 | 2000 | 2750 | 2750 | 5366 |
| Test | 2500 | 3750 | 1500 | 2761 | 2761 | 5357 |

Table 2: Dataset Distribution

models, offering a comparative analysis against the performance of each LLM under both few-shot and zero-shot prompting scenarios. The table facilitates a comprehensive evaluation by highlighting key outcomes, enabling a thorough understanding of how each model performs across the different annotation tasks and prompting methods.

## 6.1 Key Findings

Our extensive experiments revealed crucial insights, showing that Large Language Models (LLMs) are not yet fully equipped to serve as reliable annotators for the Marathi language. The disparity between LLM-based and human-generated annotations remains substantial. Even for straightforward tasks like news classification, LLM performance was suboptimal. For more complex tasks, such as the 4-class MahaHate classification, their performance was notably disappointing, as evidenced in Table 1.

Among the LLMs evaluated, GPT-4o achieved the best results compared to others, including Llama 3.1 8B, Gemma 2 (2B and 9B), and Gemini 1.0 pro. However, both open-source and closed-source LLMs exhibited notable limitations in providing accurate and reliable annotations. Our results also demonstrate that closed LLMs like GPT-4o and Gemini 1.0 Pro outperform open LLMs namely Llama and Gemma 2B and 9B but they still underperform when compared to finetuned BERT for almost all datasets.

Our evaluation ranks the models as GPT-4o > Llama 3.1 405B > Gemini 1.0 Pro > Gemma 2 9B > Llama 3.1 8B > Gemma 2 2B, highlighting that open source Llama 3.1 405B outperforms Gemini 1.0 Pro and is second only to GPT-4o.

Compared to zero-shot prompting, few-shot prompting produced more accurate results because it gave examples of the desired input-output behavior, helping the model to understand the task's context and expectations better. We observe an absolute increase in the average accuracy of few-shot prompting by 5.3%, 0.5%, 3.2%, 0.5%, 2.2%, and 0.8% compared to zero-shot prompting for models Llama 3.1 8B, Gemma 2 2B, Gemma 2 9B, Gemini 1.0 Pro, Llama 3.1 405B, and GPT-4o, respectively. While few-shot prompting techniques yielded better accuracy than zero-shot approaches, they still fell short of the performance delivered by BERT-based models.

The average accuracy gap between open-source and closed-source models is 6.7%, while the difference between closed-source models and fine-tuned BERT-based models is 13.9%, highlighting the lack of effective LLMs for low-resource languages.

BERT-based fine-tuned models on target task outperformed LLMs in the classification tasks for Marathi language because finetuning enabled better knowledge extraction and alignment with the task's requirements. On the other hand, LLMs, despite being trained on vast amounts of general data, lack performance on low resource languages and task-specific optimization, which limits their ability to extract the most relevant features for a specific task. This suggests that, despite the increasing popular-

ity of LLMs, BERT-based models continue to be highly relevant, particularly for Indic languages.

We also note that the difference in the results of BERT-based models and the LLMs is comparatively less for easy tasks like the Sentiment Analysis Task, i.e. in the MahaSent dataset. At the same time, the gap is significantly higher for complex tasks like the Hate Classification and News Classification tasks in favor of BERT-based models.

# 7 Conclusion

Our study demonstrates that while LLMs like GPT, Gemini, Gemma, and Llama show potential, they currently fall short of being reliable annotators for low-resource languages like Marathi, particularly for complex tasks. BERT-based models continue to outperform LLMs in these contexts. Moreover, these findings can be generalized to other Indic languages as well, such as Marathi, due to their morphological richness. These results indicate that further advancements are required in LLMs to make them viable alternatives to human annotations. This research highlights the need for developing more robust models tailored to the specific details of low-resource languages. This includes the creation of higher-quality, task-specific datasets for low-resource languages, ensuring better representation and reducing biases. Enhanced datasets combined with domain-specific knowledge can significantly improve annotation accuracy.

## Acknowledgement

## References

Fahd N. Al-Wesabi, Hala J. Alshahrani, Azza Elneil Osman, and Elmouez Samir Abd Elhameed. 2023. Low-resource language processing using improved deep learning with hunter–prey optimization algorithm. *Mathematics*.

Anna Bavaresco, Raffaella Bernardi, Leonardo Bertolazzi, Desmond Elliott, Raquel Fern'andez, Albert Gatt, E. Ghaleb, Mario Giulianelli, Michael Hanna, Alexander Koller, Andr'e F. T. Martins, Philipp Mondorf, Vera Neplenbroek, Sandro Pezzelle, Barbara Plank, David Schlangen, Alessandro Suglia, Aditya K Surikuchi, Ece Takmaz, and Alberto Testoni. 2024. Llms instead of human judges? a large scale empirical study across 20 nlp evaluation tasks. *ArXiv*, abs/2406.18403.

Bosheng Ding, Chengwei Qin, Linlin Liu, Lidong Bing, Shafiq R. Joty, and Boyang Albert Li. 2022. Is gpt-3 a good data annotator? In *Annual Meeting of the Association for Computational Linguistics*.

Mingqi Gao, Jie Ruan, Renliang Sun, Xunjian Yin, Shiping Yang, and Xiaojun Wan. 2023. Human-like summarization evaluation with chatgpt. *arXiv preprint arXiv:2304.02554*.

Michael A. Hedderich, Lukas Lange, Heike Adel, Jannik Strotgen, and Dietrich Klakow. 2020. A survey on recent approaches for natural language processing in low-resource scenarios. In *North American Chapter of the Association for Computational Linguistics*.

Raviraj Joshi. 2022. L3cube-mahanlp: Marathi natural language processing datasets, models, and library. *arXiv preprint arXiv:2205.14728*.

Nataliia Kholodna, Sahib Julka, Mohammad Khodadadi, Muhammed Nurullah Gumus, and Michael Granitzer. 2024. Llms in the loop: Leveraging large language model annotations for active learning in low-resource languages. *Preprint*, arXiv:2404.02261.

Panteleimon Krasadakis, Evangelos Sakkopoulos, and Vassilios S. Verykios. 2024. A survey on challenges and advances in natural language processing with a focus on legal informatics and low-resource languages. *Electronics*.

Atharva Kulkarni, Meet Mandhane, Manali Likhitkar, Gayatri Kshirsagar, and Raviraj Joshi. 2021. L3cubemahasent: A marathi tweet-based sentiment analysis dataset. In *Proceedings of the Eleventh Workshop on Computational Approaches to Subjectivity, Sentiment and Social Media Analysis*, pages 213–220.

Taja Kuzman, Igor Mozetic, and Nikola Ljubešic. 2023. Chatgpt: beginning of an end of manual linguistic data annotation. *Use Case of Automatic Genre Identification. ArXiv abs/2303.03953*.

Aishwarya Mirashi, Srushti Sonavane, Purva Lingayat, Tejas Padhiyar, and Raviraj Joshi. 2024. L3cube-indicnews: News-based short text and long document classification datasets in indic languages. *Preprint*, arXiv:2401.02254.

Saloni Mittal, Vidula Magdum, Sharayu Hiwarkhedkar, Omkar Dhekane, and Raviraj Joshi. 2023. L3cube-mahanews: News-based short text and long document classification datasets in marathi. In *International Conference on Speech and Language Technologies for Low-resource Languages*, pages 52–63. Springer.

Jay Mohta, Kenan Emir Ak, Yan Xu, and Mingwei Shen. 2023. Are large language models good annotators? In *ICBINB*.

Arbi Haza Nasution and Aytuǧ Onan. 2024. Chatgpt label: Comparing the quality of human-generated and llm- generated annotations in low-resource language nlp tasks. *IEEE Access*, 12:71876–71900.

Hrushikesh Patil, Abhishek Velankar, and Raviraj Joshi. 2022. L3cube-mahahate: A tweet-based marathi hate speech detection dataset and bert models. In *Proceedings of the Third Workshop on Threat, Aggression and Cyberbullying (TRAC 2022)*, pages 1–9.

Maja Pavlovic and Massimo Poesio. 2024. The effectiveness of llms as annotators: A comparative overview and empirical analysis of direct representation. *ArXiv*, abs/2405.01299.

Aabha Pingle, Aditya Vyawahare, Isha Joshi, Rahul Tangsali, and Raviraj Joshi. 2023. L3cube-mahasent-md: A multi-domain marathi sentiment analysis dataset and transformer models. In *Proceedings of the 37th Pacific Asia Conference on Language, Information and Computation*, pages 274–281.

Girija V R, Sudha T, and Riboy Cheriyan. 2023. Analysis of sentiments in low resource languages: Challenges and solutions. *2023 IEEE International Conference on Recent Advances in Systems Science and Engineering (RASSE)*, pages 1–6.

Aarohi Srivastava, Abhinav Rastogi, Abhishek Rao, Abu Awal Md Shoeb, Abubakar Abid, Adam Fisch, Adam R Brown, Adam Santoro, Aditya Gupta, Adrià Garriga-Alonso, et al. 2022. Beyond the imitation game: Quantifying and extrapolating the capabilities of language models. *arXiv preprint arXiv:2206.04615*.

Zhen Tan, Dawei Li, Song Wang, Alimohammad Beigi, Bohan Jiang, Amrita Bhattacharjee, Mansooreh Karami, Jundong Li, Lu Cheng, and Huan Liu. 2024. Large language models for data annotation and synthesis: A survey. *arXiv preprint arXiv:2402.13446*.

Yiming Zhu, Peixian Zhang, Ehsan-Ul Haq, Pan Hui, and Gareth Tyson. 2023. Can chatgpt reproduce human-generated labels? a study of social computing tasks. *arXiv preprint arXiv:2304.10145*.

# Speech-Based Depressive Mood Detection in the Presence of Multiple Sclerosis: A Cross-Corpus and Cross-Lingual Study

**Monica Gonzalez-Machorro**[1,2,3]**, Uwe Reichel**[1]**, Pascal Hecker** [1,4]**, Helly Hammer**[5]**,**
**Hesam Sagha**[1]**, Florian Eyben**[1,6]**, Robert Hoepner**[5]**, Björn W. Schuller**[1,2,3,7]

[1]audEERING GmbH, [2]TUM University Hospital,
[3]Munich Center for Machine Learning, [4]Hasso-Plattner Institute,
[5]Inselspital, Bern University Hospital, [6]Agile Robots, [7]Imperial College
`monica.gonzalez@tum.de`

## Abstract

Depression commonly co-occurs with neurodegenerative disorders like Multiple Sclerosis (MS), yet the potential of speech-based Artificial Intelligence for detecting depression in such contexts remains unexplored. This study examines the transferability of speech-based depression detection methods to people with MS (pwMS) through cross-corpus and cross-lingual analysis using English data from the general population and German data from pwMS. Our approach implements supervised machine learning models using: 1) conventional speech and language features commonly used in the field, 2) emotional dimensions derived from a Speech Emotion Recognition (SER) model, and 3) exploratory speech feature analysis. Despite limited data, our models detect depressive mood in pwMS with moderate generalisability, achieving a 66% Unweighted Average Recall (UAR) on a binary task. Feature selection further improved performance, boosting UAR to 74%. Our findings also highlight the relevant role emotional changes have as an indicator of depressive mood in both the general population and within PwMS. This study provides an initial exploration into generalising speech-based depression detection, even in the presence of co-occurring conditions, such as neurodegenerative diseases.

## 1 Introduction

Depression is the most common psychiatric mood disorder (World Health Organization, 2023). Its prevalence is around 5% worldwide (World Health Organization, 2023). Despite its prevalence, depression often goes untreated (Johnson et al., 2022) due to factors such as socioeconomic barriers and a shortage of healthcare professionals (Evans-Lacko et al., 2018).

Speech-based Artificial Intelligence (AI) methods offer a promising approach for fast and non-invasive screening of neurological and mental health during routine examinations (Milling et al.,

2022; Hecker et al., 2022), leveraging speech changes like reduced pitch, slower speaking rate, and articulation errors, which are common in individuals with depression (Cummins et al., 2015). These methods are accessible, scalable, and could enhance help-seeking behaviour and on-going monitoring (Johnson et al., 2022).

Prior work has utilised Machine Learning (ML) methods to detect depression using acoustic and linguistic features (Kappen et al., 2023). Mallol-Ragolta et al. (2019) trained a Recurrent Neural Network (RNN) on linguistic features for binary classification on the Distress Analysis Interview Corpus from the Wizard-of-Oz interviews (DAIC-WoZ) dataset, achieving an F1 score of 63%. Zhang et al. (2024) used wav2vec 2.0 for feature extraction and a Long Short-Term Memory (LSTM) network for binary classification using the DAIC-WoZ dataset, which yielded a 79% F1 score.

Similar work has also been conducted in other languages, such as for the German language, Menne et al. (2024) reported a balanced accuracy 88% for predicting depressive disorder against healthy controls using acoustic information, and for Italian language, in which Tao et al. (2023) reported an F1 score of 85% on the binary task of identifying depression using speech information from a reading task.

Automatic Speech Emotion Recognition (SER) research has also been effective in depression detection (Wang et al., 2020), for instance, Wang et al. (2021) developed a SER model on the DAIC-WoZ dataset for binary classification, reporting a 60% F1 score.

Depression is a common co-morbidity among people with neurodegenerative diseases, such as Multiple Sclerosis (MS), Parkinson's Disease (PD), and Alzheimer's Disease (AD), among others (Brenes, 2007), worsening both the Quality of Life (QoL) and disease prognosis (Hussain et al., 2020). In MS, for example, the lifetime risk of

depression is estimated around 50% (Arnett et al., 2008). The overlapping symptomatology of the two conditions can lead to misdiagnosis, with either one of them frequently overlooked (Hussain et al., 2020). While prior research highlights the potential of speech-based AI methods for depression detection (Cummins et al., 2015), further work is needed to assess their transferability in patients with neurodegenerative diseases like MS.

However, MS, due to its impact on the central nervous system, frequently leads to speech impairment, primarily dysarthria (Noffs et al., 2018). As a result, MS speech typically presents irregular articulatory breakdowns, distorted vowels, pitch breaks, harsh voice quality, and slow speaking rate (Noffs et al., 2018). This raises the question of whether speech-based depression detection can distinguish depressive symptoms in people with a co-existing speech impairment, such as dysarthria, due to a neurodegenerative disease, such as MS. We hypothesise that these methods would struggle to generalise and distinguish depressive symptoms in people with MS (pwMS), since some of the MS speech characteristics are similar to those found in people with depression.

This contribution aims to address this challenge by assessing the performance of common speech-based methods for depressive mood detection in pwMS. To do so, we conduct a cross-corpus and cross-lingual analysis using a well-known English-language corpus with depressive mood assessments, along with a German-language dataset of people with low MS disability, who also underwent depressive mood assessments. Our research questions are:

1. Do ML methods for depressive mood detection generalise to depressive mood detection in pwMS?

2. Given that SER models have shown promise in detecting emotional changes (Wang et al., 2021), which output from a fine-tuned SER model is more effective for depression detection: the model's final results (the classification or regression head output from a SER model) corresponding to the emotional dimensions –arousal, valence, and dominance– or the model's contextualised representations?

3. Can exploratory feature selection analysis improve generalisability of depression detection in pwMS?

This contribution is structured as follows. Section 2 introduces the datasets, features, and methods employed. Sections 3, 4, 5 present the results, limitations, and discussions. Finally, section 6 draws conclusions from the analysis.

## 2 Materials and Methods

### 2.1 Dataset

We employ two datasets: 1) The DAIC-WoZ depression dataset in English presented in (Gratch et al., 2014), and 2) a Swiss German dataset for pwMS collected under the scope of the COMMITMENT trial (Gonzalez-Machorro et al., 2023). The trial protocol was approved by national regulatory authorities and local ethic committee (BASEC-ID number 2021-02423) and registered on clinical-trials.gov (NCT05561621). The DAIC-WoZ is a collection of semi-structured interviews containing speech samples of 189 participants (Gratch et al., 2014). It provides predefined speaker-independent training, development, and testing sets, and is segmented at the turn level (Valstar et al., 2016). The dataset includes scores from the Patient Health Questionnaire-8 (PHQ-8) self-assessed depression questionnaire.

The COMMITMENT (Prediction of Non-motor Symptoms in Fully Ambulatory MS Patients Using Vocal Biomarkers) dataset consists of 50 fully ambulatory pwMS and 20 control participants. Participants with MS have low levels of disability, with a median Expanded Disability Status Scale (EDSS) score of 1.0–indicating minimal impairment– and a min/max EDSS score of 0.0/3.0, which indicates no disability to moderate disability but still walking unaided. For this paper, we only use the MS cohort. Details on the speech recordings are described in (Gonzalez-Machorro et al., 2023). Depressive mood scores for each participant are available using the Beck Depression Inventory-II (BDI-II) questionnaire. The dataset contains multiple speech tasks. However, in this paper, we utilise two spontaneous speech tasks from each patient: (1) describing the weather on the day of recording and (2) recalling a neutral memory prompted by the word "grass". These tasks are chosen because they elicit spontaneous speech and resemble the interview style of the DAIC-WoZ dataset. Data was collected using the AISoundLab web platform, which is a web app, in which each patient could navigate through a voice recording session under the supervision of a study nurse (Gonzalez-Machorro et al.,

2023). All participants provided informed consent prior to participation, and all data was pseudo-anonymised to protect patient privacy. The ethics consent unfortunately does not permit the publication of the recorded data.

In this paper, participants from the two datasets are categorised as having *depression* or *no depression* based on clinically validated threshold scores from two depression questionaries (BDI-II and PHQ-8). For the PHQ-8, participants with a score of 10 or higher are classified as having *depression* (Kroenke et al., 2001; Dhingra et al., 2011); and for the BDI-II participants with a score higher than 19 were defined as having *depression* (Beck et al., 1961). It is important to keep in mind that these scores serve as indicators of depressive symptoms rather than definitive clinical diagnoses of depression.

Audio files are downsampled to 16 kHz. Diarisation for the DAIC-WoZ data is performed using the turn-level segments provided for each speaker. A Voice Activity Recognition (VAD) algorithm[1] is applied to segment audio files from both datasets, which due to license restrictions, is not open-source. For consistency with previous work, we employ the same VAD parameter values as in (Gonzalez-Machorro et al., 2023). Transcripts are automatically obtained for each VAD segment using Whisper version 2 (Radford et al., 2023) with the *base* model for English and German language. For the DAIC-WoZ dataset, we merge the original training and development sets while the original testing set is left intact. The motivation is that due to the small dataset, we opt to use a Cross-Validation (CV) strategy for a more robust evaluation. The COMMITMENT dataset, as its purpose is purely for evaluating cross-corpus and cross-lingual generalisation, is not partitioned and it is used as an additional testing set.

Table 1 describes the metadata for both datasets across the different dataset partitions. Missing values for the questionnaires are dropped before processing. Models trained solely on the COMMITMENT dataset would likely over-fit due to insufficient participants with depressive symptoms to learn acoustic and linguistic markers of depression. Given the imbalance of the two classes, random oversampling with replacement for the two classes and a random seed of 42 is applied. To do so, we employ the package imbalanced-learn (Lemaître

---

[1]provided by audEERING GmbH

Table 1: Metadata for the two datasets employed in this study and the train-test split.

| Subset | Dataset | Total Partici-pants | Sex (F/M) | Depression / No Depression |
|--------|---------|---------------------|-----------|----------------------------|
| Train  | DAIC-WoZ | 135 | 59 / 76 | 42 / 93 |
| Test   | DAIC-WoZ COMMITMENT | 44 50 | 22 / 22 37 / 13 | 13 / 31 4 / 46 |

et al., 2017).

## 2.2 Feature extraction

We extract six commonly used acoustic and linguistic feature sets, and normalise them per dataset using the Robust Scaler, which is robust against outliers. All features are extracted at a VAD segment-level.

1. The Wav2Vec2 contextualised representations of length 1024 correspond to the mean pooling of the encoder output. These representations are extracted using a publicly available fine-tuned Wav2Vec2 model for 3-dimensional SER task (Wagner et al., 2023).

2. SER-dimensions –arousal, valence, and dominance– are obtained using the same Wav2Vec2 SER model (Wagner et al., 2023). These features represent the final outputs of the model returned by the 2-layer multitask regression head (Wagner et al., 2023). By extracting both types of information –the contextualised representations and the emotion dimensions– from the Wav2Vec2 SER model, we aim to investigate which one is more effective for depression detection.

3. Praat features (Feinberg, 2022) are extracted using Nkululeko (Burkhardt et al., 2022) and correspond to 39 features, such as voice quality, shimmer, jitter, and duration. This type of features has shown significance for depression detection (Cummins et al., 2015).

4. extended Geneva minimalistic acoustic parameter set (eGeMAPS) (Eyben et al., 2016) is extracted using the Speech & Music Interpretation by Large-space Extraction (openSMILE) feature extraction tool (Eyben et al., 2010). It contains 22 acoustic features related to prosody, voice quality, and articulation. Previous work has reported promising results in

depression detection (Cummins et al., 2015). We employ the 88 functionals and summary statistics from these features.

5. The psycholinguistic feature set consists of 51 linguistic features that represent the syntactic complexity, the proportion of sentiment tokens, and the proportion of nouns, verbs, negations, adjectives, among others.

6. RoBERTa embeddings are extracted using a multilingual model –XLM Large RoBERTa (Conneau et al., 2020) –. These embeddings correspond to the $[CLS]$ pooling output applied to the last hidden states of the model. Each segment is defined with a maximum length of 512 tokens and represented by a size of 768.

## 2.3 Methods

We define the following three modelling scenarios to investigate whether ML methods for depressive mood detection generalise in the presence of MS:

**A) Baseline Performance:** Each feature set and model type is trained and evaluated on the DAIC-WoZ training and testing sets. This task establishes a baseline for model performance in depression detection.

**B) Generalisability Evaluation:** Each feature set and model type is evaluated on the DAIC-WoZ testing set –to ensure consistent performance on the general population– and the COMMITMENT dataset. The aim is to assess how well models trained on data from the general population (DAIC-WoZ) generalise to the pwMS data.

**C) Feature Selection Modelling:** Following an exploratory feature analysis on the DAIC-WoZ training set, the resulting significant features are used for training and evaluation. This task aims to improve model performance by selecting relevant features for depression detection. Two scenarios are investigated:

**C_A)** Models are trained and evaluated on the DAIC-WoZ training and testing sets using selected features. In other words, it is Task A with selected features. This task assesses whether feature selection improves performance within the general population.

**C_B)** Models are trained on the DAIC-WoZ training set using selected features and evaluated on both the DAIC-WoZ testing set and the COMMITMENT dataset. This scenario, equivalent to Task B with selected features, explores whether feature selection improves generalisability to pwMS data.

**Exploratory feature analysis.** To investigate which features are significant to distinguish between speakers with and without depression in the training set, we use the Mann-Whitney U test ($p < 0.05$) because it is non-parametric and does not require the assumption of a normal distribution. This makes it suitable for our data, where not all features follow a normal distribution. Additionally, it is more conservative than other statistical tests, reducing the risk of Type I errors. To quantify the effect size, we use Cohen-R (Cohen, 1988). Relevant features are found by selecting among the significant ones those with an $r \geq .30$. Corrections for Type 1 errors are not performed due to the large size of the feature sets, so that the aim of this analysis is restricted to explore acoustic and linguistic feature trends.

**Modelling.** We implement supervised ML classification for implementing the three modelling tasks. For reproducibility, we seed the pseudo-random number generation. The models used are Support Vector Machine (SVM), Random Forest (RF), and eXtreme Gradient Boosting (XGB). These supervised learning algorithms were selected due to their consistently strong performance across a wide range of classification tasks (Fernández-Delgado et al., 2014). Each model is trained using Grid search 5-fold speaker-independent CV on the training set.

The hyper-parameter values optimised for the Grid Search for each model are as follows: for SVM, $C \in [10^{-4}, 10^{-3}, 10^{-2}, 10^{-1}, 1, 10]$, the kernel options include `linear` and `rbf`, and the gamma parameter is chosen from `scale` and `auto`. For XGB, the number of estimators $\in [200, 300, 450, 500]$, the learning rate $\in [0.001, 0.01, 0.1, 0.2]$, the maximum tree depth $\in [4, 5, 6]$, the column subsample ratio $\in [1, 0.3, 0.5]$, and the subsample ratio $\in [0.8, 1]$. Lastly, for the RF model, the number of estimators $\in [50, 100, 300, 500, 800, 1000]$, the criterion is either `gini` or `entropy`, the minimum number of samples required to split an internal node is $\in [2, 3]$,

and bootstrap sampling is either `True` or `False`.

The optimal hyper-parameters identified through this process are then used to train the model on the entire training set. Class weights are calculated from the training set and are incorporated to address the class imbalance in the data.

**Evaluation.** We calculate speaker-level Unweighted Average Recall (UAR), F1-score, precision, and recall. Receiver Operating Characteristic (ROC) curves and the Area Under the Curve (AUC) scores were also calculated at a speaker-level. Due to space limitations, only the ROC curves for the best-performing tasks are presented. We also compute the 95% Confidence Interval (CI) for the UAR. The CIs were calculated using 1000 bootstrapping iterations [2].

## 3 Results

### 3.1 Exploratory feature analysis

The Mann-Whitney U test is applied to each feature in the training set of the DAIC-WoZ dataset. Due to interpretability limitations, the Wav2Vec2 and the RoBERTa representations are excluded from the analysis. The number of significant ($p < 0.05$) features with a sufficiently high effect size ($r \geq 0.30$) identified per feature set are: 1) SER-dimensions: 1 feature–valence–; 2) Praat features: 33 out of 39 features; 3) eGeMAPS: 64 out of 88 functionals; 4) Psycholinguistic feature set: 18 out of 51 features. These selected features are used in the modelling task C_A and C_B to assess whether feature selection improves modelling performance. Figure 1 shows the valence distributions for the binary depression class ("no_depression" and "depression"), which is the only significant features found for the SER-dimensions.

### 3.2 Modelling Results

Table 2 shows UAR and its CIs, F1-score, precision, and recall for *depression* (Dep.) and *no depression* (No Dep.) classes, across the best-performing models and all feature sets. As we are tackling a binary classification problem, the chance-level UAR is 50%. The best result for Task A (Baseline Performance) with acoustic features is achieved using SVM and SER-dimensions (UAR: 73%), while the best result with linguistic features is achieved using SVM and RoBERTa embeddings (UAR: 56%). For

Figure 1: Feature distributions for the binary depression class depression class and valence dimension from SER in the DAIC-WoZ training set. This feature presents a $r$ of 0.66 and $p < 0.001$.

Task B (Generalisability Evaluation), Wav2Vec2 embeddings and Psycholinguistic features achieved the best performances (UAR: 66% and 62%, respectively). SER-dimensions in Task B show a performance drop. For Tasks C_A and C_B (Feature Selection Modelling on Tasks A and B), XGB with SER-dimensions obtained the highest UARs of 79% and 74%, respectively. Since SER-dimensions shows consistently good performance in all tasks, Figure 2 shows the ROC curves and AUC values for all tasks.

## 4 Discussion

In this paper, we explore three research questions: **1)** Do ML methods for depressive mood detection generalise to depressive mood detection in pwMS? Results in Table 2 indicate that for Task B (Generalisability Evaluation), acoustic-based features show reasonable generalisability to distinguish depression in pwMS, with only a modest performance decline compared to results from Task A (Baseline Performance).

In the case of the Wav2Vec2 features, a drop in

Table 2: Speaker-level test results. **A**: Baseline Performance. **B**: Generalisability Evaluation. **C_A**: Feature Selection on Task A. **C_B**: Feature Selection on Task B. The best-performing combinations for acoustic-based models are marked in **bold** and *; and linguistic models as **bold**[†]. Dep. corresponds to the *depression* class and No Dep. correspond to *no depression*.

| Task | Feature | Model | UAR[%] | F1[%] | Precision[%] | | Recall[%] | |
|------|---------|-------|--------|-------|------|------|------|------|
| | | | | | Dep. | No Dep. | Dep. | No Dep. |
| A | Wav2Vec2 | XGB | 66(49-81) | 65 | 81 | 47 | 71 | 62 |
| | SER-dimensions | SVM | **73(57-84)*** | **67*** | 48 | 90 | 85 | 61 |
| | Praat | XGB | 49(42-62) | 45 | 70 | 25 | 90 | 8 |
| | eGeMAPS | XGB | 54(46-69) | 53 | 72 | 50 | 94 | 15 |
| | Psycholinguistic | SVM | 46(32-63) | 46 | 68 | 25 | 61 | 31 |
| | RoBERTa | SVM | **56(48-71)**[†] | **54**[†] | 67 | 73 | 15 | 97 |
| B | Wav2Vec2 | SVM | **66(54-80)*** | **67*** | 50 | 88 | 41 | 91 |
| | SER-dimensions | SVM | 64(50-76) | 57 | 28 | 89 | 65 | 64 |
| | Praat | SVM | 47(33-60) | 39 | 16 | 79 | 53 | 40 |
| | eGeMAPS | XGB | 56(49-69) | 57 | 43 | 84 | 18 | 95 |
| | Psycholinguistic | SVM | **62(48-74)**[†] | **54**[†] | 26 | 88 | 65 | 60 |
| | RoBERTa | SVM | 55(49-67) | 54 | 50 | 83 | 12 | 97 |
| C_A | SER-dimensions | XGB | **79(70-87)*** | **70*** | 50 | 100 | 100 | 58 |
| | Praat | SVM | 51(36-68) | 51 | 31 | 71 | 31 | 71 |
| | eGeMAPS | XGB | 58(48-74) | 58 | 60 | 74 | 23 | 94 |
| | Psycholinguistic | SVM | 48(33-66) | 48 | 69 | 28 | 58 | 38 |
| C_B | SER-dimensions | XGB | **74(60-84)*** | **65*** | 37 | 93 | 76 | 71 |
| | Praat | SVM | 46(32-59) | 38 | 16 | 79 | 53 | 39 |
| | eGeMAPS | XGB | 57(46-70) | 56 | 29 | 84 | 29 | 84 |
| | Psycholinguistic | SVM | 54(42-68) | 51 | 22 | 84 | 41 | 68 |



Figure 2: ROC curve and AUC value at a speaker-level for the best-performing models using the SER-dimenions as feature set across all tasks. Task A: Baseline Performance. B: Generalisability Evaluation. C_A: Feature Selection on Task A. C_B: Feature Selection on Task B.

performance for the two tasks is not found, which suggests that these features are transferable to other languages and groups with other co-morbidities such as MS. Interestingly, in the case of the eGeMAPS features, a minimal increase in performance is observed in Task B, which also suggests a generalisability capacity.

For Tasks C_A and C_B (Feature Selection Modelling), similar patterns are observed as in Tasks A and B, with SER-dimensions consistently outperforming other feature sets and demonstrating strong transferability in detecting depression among pwMS. This is further illustrated in Figure 2, which highlights the effectiveness of SER-dimensions in the context of MS.

The top-performing results for Tasks A (using SER-dimensions) and B (using Wav2Vec2 features) demonstrate greater precision in predicting the absence of depression (90% for "No Dep." in Task

A; 88% for "No Dep." in Task B) compared to predicting depression. This finding indicates that identifying depression using speech presents similar challenges in both same-language and cross-lingual contexts, as well as in the general population and among groups with co-morbidities, such as MS.

Interestingly, RF models did not outperform XGB or SVM in any task or feature set; consequently, they are excluded from Table 2. This was already reported by (Fernández-Delgado et al., 2014), where XGB has been shown to outperform RF in many cases.

**2)** Given that SER models have shown promise in detecting emotional changes, which output from a fine-tuned SER model is more effective for depression detection: the model's final predictions corresponding to the emotional dimensions or the model's contextualised representations? As shown in Table 2, the SER-dimensions and Wav2Vec2 representations achieve the highest UAR for Task A and Task B, respectively. SER-dimensions also outperform all other feature sets in Task C reaching the highest performance. Likely due to the high dimensionality of the Wav2Vec2 embeddings, SER-dimensions show overall better results by a small margin. However, the performance of SER-dimensions and Wav2Vec2 features heavily relies on the performance of the underlying SER model (Wagner et al., 2023), which was finetuned using the MSP-Podcast dataset (English language) (Lotfian and Busso, 2019). It is, therefore, unclear the cross-lingual generalisability of these features when training data would include languages other than English.

**3)** Can feature selection improve generalisability of depression detection in pwMS? Results for acoustic-feature-based models, with the exception of the Praat features, suggest that indeed, feature selection can improve the performance of depression detection. The feature analysis for SER-dimensions reveals that only valence among the three dimensions is significantly predictive, highlighting its important role as an indicator of depression in both the general population and pwMS. This finding is illustrated in Figure 2, which shows that individuals without depressive symptoms tend to use higher positive valence in spontaneous interviews compared to those with depressive symptoms. This aligns with prior research, such as (Trifu et al., 2024), which found that individuals with de-

pression display lower positive valence than those without. This pattern may be attributed to a core symptom of depression: emotional dysfunction characterised by a predominant negative emotional state (Yang et al., 2023).

## 5 Limitations

In the case of text-based models, RoBERTa embeddings achieve above-chance performance in both Task A and Task B while psycholinguistic-feature-based models exhibit an unexpected trend: their performance on Task B surpassed that of Task A, C_A, and C_B. The suboptimal performances of text-based models may be due to the use of VAD segments for feature extraction, which ensured a consistent preprocessing pipeline across acoustic and text features, enabling direct comparisons between model types in detecting depression. While VAD segments effectively captured acoustic cues, contributing to strong performances, their short duration may have been less optimal for text-based features, such as word class proportions, which benefit from longer discourse contexts. The language-specific nature of these features also might have contributed to their struggle to generalise to the German-speaking MS population. Future work should explore longer segments to optimise text-based models, building on this study's foundation.

A limitation of this contribution arises from the use of different languages, recording conditions, and depression assessments. Although we try to tackle this by feature normalisation and the restriction to spontaneous speech, further research should explore the impact of language, depression assessments, and recording variations on the generalisability of speech-based depression detection. In this paper, we cannot definitively differentiate the extent to which the drop in model performance when evaluated on the MS population is influenced by language differences, recording conditions or the presence of MS itself.

Moreover, since both MS and depression are heterogenous conditions (Gaitán and Correale, 2019), implementing personalised approaches when screening for depression in pwMS is a crucial next step. Future work should also explore different stages of MS –this study focused on low-disability patients– and account for other co-morbidities in MS, like fatigue and cognitive decline, which may also influence speech. Also, the MS cohort was receiving pharmacological treatment, including com-

mon antidepressants for those MS patients diagnosed with depression, that could influence mood and, consequently, speech patterns. Although the general population diagnosed with depression from the DAIC-WoZ dataset may also have been undergoing pharmacological treatments, this information is not available in the dataset, preventing analysis of this potential confounding factor.

To further evaluate the transferability of speech-based depression detection, it is important to examine other common diseases where depression is a common co-morbidity and speech is impacted, such as PD or AD. A lack of depression scores in speech datasets for these disorders is a major limitation in this regard. Finally, acoustic and linguistic features alone cannot fully capture the multifaceted nature of depression. These ML methods are intended to augment established screening approaches. Incorporating other bio-signals, such as physiological data, could not only enhance performance but also provide a more comprehensive understanding of the disorder.

## 6 Conclusion

In this cross-corpus and cross-lingual study, we explore the efficacy of speech-based depressive mood detection in the presence of MS and across English and German languages. Our findings highlight the significance of emotional dimensions –arousal, valence, and dominance– in identifying depressive symptoms, not only in the general population but also within pwMS. Additionally, acoustic feature sets like eGeMAPS also demonstrate potential for generalisability in this context. However, further research is needed to establish robust conclusions. This study, despite its limitations, represents a step forward towards the integration and generalisability of speech-based depression detection methods. Non-invasive speech-based AI systems for depression detection hold the potential to improve the QoL for individuals with this disorder, even in the presence of other illnesses.

## Ethics Statement

## References

Peter A. Arnett, Fiona H. Barwich, and Joe E. Beeney. 2008. Depression in multiple sclerosis: Review and theoretical proposal. *Journal of the International Neuropsychological Society*, 14(5):691–724.

Aaron T. Beck, Clyde H. Ward, Myer Mendelson, John Mock, and John Erbaugh. 1961. An inventory for measuring depression. *Archives of General Psychiatry*, 4(6):561–571.

Gretchen A Brenes. 2007. Anxiety, depression, and quality of life in primary care patients. *Primary care companion to the Journal of clinical psychiatry*, 9(6):437–443.

Felix Burkhardt, Johannes Wagner, Hagen Wierstorf, Florian Eyben, and Björn Schuller. 2022. Nkululeko: A tool for rapid speaker characteristics detection. In *2022 Language Resources and Evaluation Conference, LREC 2022*, pages 1925–1932. European Language Resources Association (ELRA).

Jacob Cohen. 1988. *Statistical Power Analysis for the Behavioral Sciences*, 2nd edition. Erlbaum, Hillsdale, NJ.

Alexis Conneau, Kartikay Khandelwal, Naman Goyal, Vishrav Chaudhary, Guillaume Wenzek, Francisco Guzmán, Edouard Grave, Myle Ott, Luke Zettlemoyer, and Veselin Stoyanov. 2020. Unsupervised cross-lingual representation learning at scale. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 8440–8451, Online. Association for Computational Linguistics.

Nicholas Cummins, Stefan Scherer, Jarek Krajewski, Sebastian Schnieder, Julien Epps, and Thomas F. Quatieri. 2015. A review of depression and suicide risk assessment using speech analysis. *Speech Communication*, 71:10–49.

Subash S. Dhingra, Kurt Kroenke, Matthew M. Zack, Tara W. Strine, and Lina S. Balluz. 2011. PHQ-8 Days: A Measurement Option for DSM-5 Major Depressive Disorder (MDD) Severity. *Population Health Metrics*, 9:11.

Sara Evans-Lacko, Sergio Aguilar-Gaxiola, Ahmad Al-Hamzawi, Jordi Alonso, Corina Benjet, Ronny Bruffaerts, Wai Tat Chiu, Silvia Florescu, Giovanni de Girolamo, Oye Gureje, Josep Maria Haro, Yanling He, Chiyi Hu, Elie G. Karam, Norito Kawakami, Sing Lee, Crick Lund, Viviane Kovess-Masfety, Daphna Levinson, Fernando Navarro-Mateu, and Graham Thornicroft. 2018. Socio-economic variations in the mental health treatment gap for people with anxiety, mood, and substance use disorders: results from the who world mental health (wmh) surveys. *Psychological Medicine*, 48(9):1560–1571.

Florian Eyben, Klaus R. Scherer, Björn W. Schuller, Johan Sundberg, Elisabeth André, Carlos Busso, Laurence Y. Devillers, Julien Epps, Petri Laukka,

Shrikanth S. Narayanan, and Khiet P. Truong. 2016. The geneva minimalistic acoustic parameter set (gemaps) for voice research and affective computing. *IEEE Transactions on Affective Computing*, 7(2):190–202.

Florian Eyben, Martin Wöllmer, and Björn Schuller. 2010. opensmile - the munich versatile and fast open-source audio feature extractor. In *Proceedings of the 18th ACM International Conference on Multimedia (ACM MM)*, pages 1459–1462, Florence, Italy. ACM.

David R Feinberg. 2022. Parselmouth praat scripts in python.

Manuel Fernández-Delgado, Eva Cernadas, Senén Barro, and Dinani Amorim. 2014. Do we need hundreds of classifiers to solve real world classification problems? *The journal of machine learning research*, 15(1):3133–3181.

María I. Gaitán and Jorge Correale. 2019. Multiple sclerosis misdiagnosis: A persistent problem to solve. *Frontiers in Neurology*, 10.

Monica Gonzalez-Machorro, Pascal Hecker, Uwe D. Reichel, Helly N. Hammer, Robert Hoepner, Lisa Pedrotti, Alisha Zmutt, Hesam Sagha, Johan van Beek, Florian Eyben, Dagmar M. Schuller, Björn W. Schuller, and Bert Arnrich. 2023. Towards Supporting an Early Diagnosis of Multiple Sclerosis using Vocal Features. In *Proc. INTERSPEECH 2023*, pages 1518–1522.

Jonathan Gratch, Ron Artstein, Gale Lucas, Giota Stratou, Stefan Scherer, Angela Nazarian, Rachel Wood, Jill Boberg, David DeVault, Stacy Marsella, David Traum, Skip Rizzo, and Louis-Philippe Morency. 2014. The distress analysis interview corpus of human and computer interviews. In *Proceedings of the Language Resources and Evaluation Conference (LREC)*.

Pascal Hecker, Nico Steckhan, Florian Eyben, Björn W. Schuller, and Bert Arnrich. 2022. Voice analysis for neurological disorder recognition–a systematic review and perspective on emerging trends. *Front. Digit. Health*, 4:842301.

Madiha Hussain, Prabhat Kumar, Sara Khan, Domonick K Gordon, and Safeera Khan. 2020. Similarities between depression and neurodegenerative diseases: Pathophysiology, challenges in diagnosis and treatment options. *Cureus*, 12(11):e11613.

Jemimah A. Johnson, Prachi Sanghvi, and Seema Mehrotra. 2022. Technology-based interventions to improve help-seeking for mental health concerns: A systematic review. *Indian J. Psychol. Med.*, 44(4):332–340.

Mitchel Kappen, Marie-Anne. Vanderhasselt, and George M. Slavich. 2023. Speech as a promising biosignal in precision psychiatry. *Neuroscience Biobehavioral Reviews*, 148:105121.

Kurt Kroenke, Robert L. Spitzer, and Janet B. W. Williams. 2001. The phq-9: Validity of a brief depression severity measure. *Journal of General Internal Medicine*, 16(9):606–613.

Guillaume Lemaître, Fernando Nogueira, and Christos K. Aridas. 2017. Imbalanced-learn: A python toolbox to tackle the curse of imbalanced datasets in machine learning. *Journal of Machine Learning Research*, 18(17):1–5.

Reza Lotfian and Carlos Busso. 2019. Building naturalistic emotionally balanced speech corpus by retrieving emotional speech from existing podcast recordings. *IEEE Transactions on Affective Computing*, 10(4):471–483.

Adria Mallol-Ragolta, Ziping Zhao, Lukas Stappen, Nicholas Cummins, and Björn W. Schuller. 2019. A Hierarchical Attention Network-Based Approach for Depression Detection from Transcribed Clinical Interviews. In *Proceedings of Interspeech 2019*.

Felix Menne, Felix Dörr, Julia Schräder, Johannes Tröger, Alexandra König, and Lisa Wagel. 2024. The voice of depression: speech features as biomarkers for major depressive disorder. *BMC Psychiatry*, 24:794.

Manuel Milling, Florian B. Pokorny, Katrin D. Bartl-Pokorny, and Björn W. Schuller. 2022. Is speech the new blood? recent progress in ai-based disease detection from audio in a nutshell. *Frontiers in Digital Health*, 4.

Gustavo Noffs, Thushara Perera, Scott C. Kolbe, Camille J. Shanahan, Frederique M.C. Boonstra, Andrew Evans, Helmut Butzkueven, Anneke van der Walt, and Adam P. Vogel. 2018. What speech can tell us: A systematic review of dysarthria characteristics in multiple sclerosis. *Autoimmunity Reviews*, 17(12):1202–1209.

Alec Radford, Jong Wook Kim, Tao Xu, Greg Brockman, Christine McLeavey, and Ilya Sutskever. 2023. Robust speech recognition via large-scale weak supervision. In *Proceedings of the 40th International Conference on Machine Learning*, ICML'23. JMLR.org.

Fuxiang Tao, Anna Esposito, and Alessandro Vinciarelli. 2023. The androids corpus: A new publicly available benchmark for speech based depression detection. In *Proc. Interspeech*, pages 4149–4153.

Raluca Nicoleta Trifu, Bogdan Nemeș, Dana Cristina Herta, Carolina Bodea-Hategan, Dorina Anca Talaș, and Horia Coman. 2024. Linguistic markers for major depressive disorder: a cross-sectional study using an automated procedure. *Frontiers in Psychology*, Volume 15 - 2024.

Michel Valstar, Maja Pantic, Jonathan Gratch, Björn Schuller, Fabien Ringeval, Denis Lalanne, Mercedes Torres Torres, Stefan Scherer, Giota Stratou, and Roddy Cowie. 2016. Avec 2016: Depression, mood, and emotion recognition workshop and challenge. In

*Proceedings of the 6th international workshop on audio/visual emotion challenge*, pages 3–10.

Johannes Wagner, Andreas Triantafyllopoulos, Hagen Wierstorf, Maximilian Schmitt, Felix Burkhardt, Florian Eyben, and Björn W Schuller. 2023. Dawn of the transformer era in speech emotion recognition: Closing the valence gap. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, pages 1–13.

Hongbo Wang, Yu Liu, Xiaoxiao Zhen, and Xuyan Tu. 2021. Depression speech recognition with a three-dimensional convolutional network. *Frontiers in Human Neuroscience*, 15.

Xusheng Wang, Xing Chen, and Congjun Cao. 2020. Human emotion recognition by optimally fusing facial expression and speech feature. *Signal Processing: Image Communication*, 84:115831.

World Health Organization. 2023. Depression. https://www.who.int/news-room/fact-sheets/detail/depression. Accessed: 2025-03-10.

Chaoqing Yang, Xinying Zhang, Yuxuan Chen, Yunge Li, Shu Yu, Bingmei Zhao, Tao Wang, Lizhu Luo, and Shan Gao. 2023. Emotion-dependent language featuring depression. *Journal of Behavior Therapy and Experimental Psychiatry*, 81:101883.

Xu Zhang, Xiangcheng Zhang, Weisi Chen, Chenlong Li, and Chengyuan Yu. 2024. Improving speech depression detection using transfer learning with wav2vec 2.0 in low-resource environments. *Scientific Reports*, 14:9543.

# Signals from Academic Hiring: A Decade of Skill Demands at a Danish University

**Zhiru Sun**
University of Southern Denmark
Department of Design, Media, and Educational Science
zhiru@sdu.dk

**Jakob Mørup Wang**
AI Consultant
Aarhus Kommune
jakob@skipconnections.com

## Abstract

The rapid emergence of Generative AI (GenAI) is transforming labor market expectations and prompting a re-evaluation of skill priorities in higher education. This study investigates how academic skill demands have evolved over the past decade by analyzing job postings from a major Danish university (2013–2023) using natural language processing (NLP) and statistical modeling. Leveraging the ESCO taxonomy, we extract and classify skills into digital, research, and transversal categories. Our findings reveal a recent shift in hiring strategy towards fewer but more skill-intensive roles, a declining emphasis on digital skills, and a rise in research-oriented skills. Additionally, we observe significant disciplinary variation, with Engineering emphasizing digital skills, while Social Sciences prioritizing research competences. In the Humanities, emerging skills increasingly reflect demands in societal engagement and digital literacy. These results offer data-driven insights into the alignment of curricula with evolving labor market demands in the GenAI era.

## 1 Introduction

The rapid rise of Generative Artificial Intelligence (GenAI), exemplified by the launch of ChatGPT in late 2022, has profoundly impacted educational landscapes and the global labor market (Adiguzel et al., 2023; Johnson et al., 2021). As AI tools increasingly assist – or even automate – complex intellectual tasks, educators around the world are grappling with fundamental questions: *What skills should we teach to help students thrive in an AI-driven world?* Much of the current educational research has focused on supply-side adaption. A growing body of research has underscored the importance of cultivating AI literacy across disciplines (Ng et al., 2021), highlighting skills such as prompt engineering (Walter, 2024), computational thinking (Weng et al., 2024; Dohn et al., 2022),

critical thinking(Muthmainnah et al., 2022). While these contributions are vital, one crucial question remains underexplored. *What does the labor market expect from future graduates in the age of AI?* This study addresses the question from the perspective of employer demand. We examined academic job postings from a major Danish university between 2013 and 2023 using natural language processing (NLP) and statistical modeling. By tracing how skill requirements have evolved over time and across disciplines, our aim is to provide evidence-based guidance on the skills that academic institutions may need to prioritize to prepare students for the GenAI era.

## 2 Literature Review

Globalization and digitization are reshaping labor markets, with AI emerging as a key driver for workforce transformation (Johnson et al., 2021). According to *European Center for the Development of Vocational Training* (Cedefop, 2023), nearly half of European adult workers encountered new digital technologies in their workplaces during 2020-2021, with 35% needing to upskill. The *World Economic Forum* (2025) similarly projects that 39% of core skills will change by 2030, driven by advances in automation, AI integration, and digital transformation.

In this evolving landscape, demand is rising for both *technical skills* – such as AI and big data, and digital literacy – and *transversal skills* including resilience, flexibility, and leadership. Research in labor economics reinforces this dual emphasis in the AI-driven labor market. On the technical side, skills in machine learning, natural language processing, computer vision, big data analytics, and programming are widely recognized as essential for AI development and application (Alekseeva et al., 2021; Johnson et al., 2021). Equally important, employers are also increasingly valuing transversal skills such as adaptability, resilience, and creativity

(Babashahi et al., 2024; Asylbekova et al., 2023). These human-centric skills are considered more resistant to automation and are critical for navigating dynamic, AI-infused work environments (Poláková et al., 2023; Belchior-Rocha et al., 2022).

Despite these findings, a recurring concern in the literature is the growing mismatch between skill supply and demand (Basson et al., 2023; Singh Dubey et al., 2022; Pater et al., 2022). Employers across sectors face significant challenges in recruiting employees with the right combination of technical and transversal skills to harness the full potential of AI (Sidhu et al., 2024). In addition, most existing studies focus primarily on the private sector of information and communication technology (ICT) and rely heavily on data from the US labor market (Babashahi et al., 2024; Johnson et al., 2021), limiting their relevance to other sectors and national contexts.

This study addresses these gaps by shifting focus to the public sector - specifically academic employment in Denmark, a highly digitalized country with robust labor market data. Academic job markets are of particular interest because they encompass a wide range of disciplines, from engineering and health sciences to the humanities and social sciences. Moreover, they serve a dual function: both responding to labor market demands and influencing future skill supply through hiring decisions, curriculum development, and research priorities.

## 3 Research Aim and Questions

This study aims to examine how skill demands in academic employment have evolved over the past decade in response to technological change. We analyze academic job advertisements from a major Danish university between 2013 to 2023. Using the European Skills, Competences, Qualifications and Occupations (ESCO) taxonomy, we extract specific skills mentioned in each job ad and classify them into three broad skill types:

- **Digital skills** – Technical and computational competences, such as programming, data analysis, and proficiency with digital tools;

- **Research skills** – Analytical, investigative, and discipline-specific competencies central to academic inquiry;

- **Transversal skills** – Interpersonal and cognitive skills such as communication, teamwork, adaptability, and problem-solving.

By tracking changes in both individual skills and skill types over time and across academic disciplines, the study offers empirical insights into how skill demands in academia have evolved. It also offers data-driven guidance to align curricula with emerging labor market needs. The study was guided by the following research questions:

1. How have the volume of academic job ads and the average number of skill requirements per ad changed over time? What is the relationship between these two trends?

2. How have demands for digital, research, and transversal skills evolved over the past decade in academic hiring?

3. Are there significant differences in mean skill demands across faculties and over time, both within and across the three skill types?

4. Which specific skills are growing and declining in demand over time, both across all faculties and within the humanities in particular?

## 4 Methodology

### 4.1 NLP: Skill Extraction and Classification

We developed an end-to-end pipeline to extract and classify skills from unstructured job advertisements based on ESCO taxonomy. The pipeline involved six stages: (1) data collection and pre-processing, (2) taxonomy alignment and encoder selection, (3) retrieval-augmented supervision with a large language model (LLM), (4) multi-objective fine-tuning of a sentence encoder, (5) normalized skill distribution at the job-ad level, and (6) evaluation. An overview is shown in Figure 1.

The pipeline began with web scraping of over 3 million Danish job ads from Jobindex, Denmark's largest job portal. We processed the data to isolate skill-relevant content through HTML parsing, regular expressions, and sentence segmentation. For each sentence, we generated semantic embeddings using the multilingual sentence encoder *paraphrase-multilingual-mpnet-base-v2*. We then retrieved the 25 most similar ESCO skills based on cosine similarity. These candidate skills were embedded into structured prompts used to query GPT-4o mini, which determined whether the sentence expresses a skill requirement and, if so, which ESCO skills are relevant. This retrieval-augmented prompting procedure resulted in more

than 536,767 labeled sentences, covering 76.7% of ESCO's 13,896 skills. Our approach represents a novel contribution to the NLP community by combining dense sentence embeddings, retrieval-based skill alignment, and LLM-based labeling in a semi-supervised pipeline for domain-specific skill extraction.

To improve sentence representations for skill recognition, we fine-tuned the sentence encoder using a multi-objective loss function: (1) a binary classification task to predict whether a sentence contains a skill, and (2) a ranking loss to align sentence embeddings with relevant ESCO skill embeddings. Fine-tuning incorporated stratified sampling, layer-wise learning rate decay, and early stopping based on ranking performance metrics.

Finally, each job ad was represented with a normalized probability distribution over ESCO skills. This probabilistic representation captured the relative importance of multiple skills within each ad while accounting for differences in length and verbosity. Full technical details of the pipeline's development and implementation are available in the associated methodological paper (Authors, 2025) .

## 4.2 Data application: Academic Job Postings in Denmark

We applied the pipeline to academic job advertisements from a major Danish university from 2013 to 2023. For each job ad, the pipeline extracted a list of relevant ESCO skills, each assigned a probability score that reflects its importance within the job ad.

To contextualize skill demand, we further enriched the dataset by identifying both faculty affiliation and position type. Faculty information was extracted by scanning the job descriptions for predefined keywords (e.g., "faculty of humanities", "humanistiske fakultet"). Through a keyword-matching function, each job was assigned to one of the five faculties or labeled as *unknown* if no match was detected. A similar rule-based approach was used to determine position types (e.g., "assistant professor", "postdoc", "PhD") based on keyword detection.

Based on ESCO's taxonomy, we then assigned each extracted skill to one of three predefined skill types: *digital, research, or transversal*. Skills falling outside these categories were labeled as *other*. This structured classification enabled us to conduct statistical analyses of skill demand evolution over time, and across faculty and academic position.

## 4.3 Statistical Analysis

To address RQ3, we used a linear mixed-effects model (LMM) to assess how average skill demands per ad have changed over time across faculties and skill types. LMMs are well-suited for hierarchical data, as they accommodate both fixed and random effects and allow for repeated observations within groups (Gelman and Hill, 2021). The dependent variable was the mean number of skills per job ad. Fixed effects included year, faculty, and skill type, along with their two-way interactions (i.e., year × faculty, year × skill type, faculty × skill type). Random intercepts were specified for each faculty × skill type combination. Model estimation used restricted maximum likelihood (REML) and Satterthwaite's approximation for degrees of freedom.

To address RQ4 on identifying growing and declining skills, we analyzed longitudinal trends in the relative frequency of skill mentions. For each year from 2013 to 2023, we calculated the proportion of job ads mentioning each skill. We then used ordinary least squares (OLS) linear regression to estimate the slope of change over time. A positive slope indicates increasing demand (growing skill), while a negative slope signals decreasing demand (declining skill). This trend analysis approach aligns with established practices in time-series modeling (Montgomery et al., 2021).

## 5 Results

## 5.1 RQ1: Job volume and skill intensity over time

Figure 2 illustrates notable fluctuations in job volume and average number of skills per job ad from 2013 to 2023. Between 2016 to 2020, number of job postings increased steadily, peaking in 2019. Meanwhile, the average skill demand per ad dropped sharply in 2016 and gradually rose again through 2019. This pattern suggests that, during periods of strong hiring demand, employers may have relaxed skill requirements to attract a broader application pool.

In 2021, job postings declined sharply, largely due to the COVID-19 pandemic. However, the average number of skills per ad remained stable and even slightly increased, indicating a shift toward fewer but more skill-intensive job positions. From 2022, job volume began to slowly recover, while av-

erage skill requirements per ad rose sharply, reaching their peak in 2023. This divergent pattern implies a transition toward more selective and skill-intensive hiring practices.

## 5.2 RQ2: Evolution of skill type

Figure 3 presents trends in the average number of skills per ad, categorized by ESCO's three skill types: *digital, research, and transversal*. Between 2013 to 2019, digital skills were the most frequently mentioned skill type in job ads, peaking in 2019. However, their prominence declined sharply thereafter. Meanwhile, research skills remained relatively stable until 2020, then increased steadily and surpassed digital skills by 2022. Transversal skills remained consistently low throughout the period, with only a modest increase after 2021. These trends suggest a shift in hiring priorities from digital proficiency toward research-oriented expertise.

## 5.3 RQ3: Differences by faculty and skill type over time

LLM results revealed significant variation in skill demands by faculty and skill type, with important interaction effects over time. While the main effect of year was not statistically significant ($p=0.53$), interactions between year and faculty ($p=0.49$), and between year and skill type ($p=0.22$) were also not statistically significant, suggesting limited evidence that temporal trends differ across disciplines or skill types.

Across the full-time span (2013–2023), digital skills were the most emphasized ($M = 3.02$), followed by research skills ($M = 2.28$), and then transversal skills ($M = 0.25$). All pairwise differences between skill types were statistically significant ($p < .001$, Bonferroni-adjusted). The faculty × skill type interaction was also highly significant ($p < .001$), showing distinct disciplinary profiles: Engineering emphasized digital skills most strongly, while Social Sciences prioritized research skills. The Humanities and Natural Sciences exhibited relatively balanced, but lower overall skill intensities. Detailed trends are shown in Figure 4.

## 5.4 RQ4: Growing and declining skills

To identify long-term trends in specific skill demands, we ranked all extracted skills by the slope of their linear trend from 2013 to 2023. The top 10 skills with the most positive and most negative slopes were classified as growing and declining skills, respectively. This analysis was conducted for both the full dataset and a subset focused on humanities faculty (Figure 5 and Figure 6).

In the humanities, emerging skills included "apply knowledge of social sciences and humanities," "web analytics," and "media studies"—all of which reflect a growing demand for interdisciplinary, data-informed, and applied research competencies. Conversely, declining skills included more traditional academic and administrative tasks, such as "assist students with their enrolment" and "contribute to specialised publications". These patterns suggest a reconfiguration of academic roles toward greater societal engagement and digital literacy.

## 6  Discussion

This study examined how skill requirements in academic job postings have evolved over the past decade in response to technological change. The findings offer several key insights relevant to curriculum development and institutional strategy.

The analysis of job volume and skill intensity reveals a shift in hiring strategy. While the number of academic positions has declined since the pandemic, the average number of skills required per job has increased, particularly after 2022. This suggests a move toward more selective recruitment, with greater emphasis on multi-skilled candidates.

Trends in skill types indicate a significant change in hiring priorities. Digital skills, once dominated, have declined in emphasis since 2020, while research skills becoming more prominent. This shift may reflect an institutional assumption that basic digital literacy is now a baseline expectation, with greater value placed on disciplinary depth and research capability.

Faculty-level analysis highlights the need for discipline-specific strategies. Engineering continues to prioritize digital skills, while Social Sciences emphasize research competencies. The Humanities and Natural Sciences show lower overall skill intensities, with more balanced distributions. At the same time, the identification of growing and declining skills provides concrete evidence of how academic expectations are shifting.

While this study draws on job postings from only one major Danish university, such postings serve as a clear demand-side signal that can guide curriculum review. They make explicit the competencies institutions prioritize in recruitment—whether discipline-specific research expertise, technical skills, or transversal abilities. These signals can

support curriculum decision-making in several ways: (1) highlighting skills that are in demand but underrepresented in existing courses; (2) identifying emerging competencies that could be incorporated into electives or interdisciplinary modules; and (3) providing evidence to inform program evaluation, accreditation, and strategic planning.

These findings also underscore the importance of systematically tracking emerging skills—especially in the current AI-driven era, where technological capabilities and work practices evolve at unprecedented speed. University curricula, constrained by structural factors such as accreditation cycles, lengthy program approval processes, and institutional governance, often lag behind the pace at which AI-related competencies emerge. Longitudinal analysis of skill trends, such as that presented here, can serve as a strategic foresight tool: rather than prompting reactive changes to every short-term fluctuation, it can help identify persistent, multi-year patterns that signal more durable shifts in competence demand. By distinguishing between fleeting spikes and sustained trends, institutions can prioritize curriculum updates that are both timely and resilient, equipping students with adaptive, interdisciplinary, and AI-literate skill sets—preparing them not only to navigate but to actively contribute to the rapidly evolving landscape of the GenAI era.

Taken together, these results contribute to the ongoing debate around what skills higher education should prioritize in the GenAI era. By offering a demand-side, data-driven view of how academic skill requirements have evolved, this study provides actionable insights for aligning university curricula with labor market expectations and making educational systems more adaptive, interdisciplinary, and future-oriented.

# References

Tufan Adiguzel, Mehmet Haldun Kaya, and Fatih Kürsat Cansu. 2023. Revolutionizing education with ai: Exploring the transformative potential of chatgpt. *Contemporary Educational Technology*, 15(3).

Liudmila Alekseeva, José Azar, Mireia Giné, Sampsa Samila, and Bledi Taska. 2021. The demand for ai skills in the labor market. *Labour economics*, 71:102002.

MP Asylbekova, TN Otarova, and DC Yelkin. 2023. The importance of transversal skills in higher education curricula and in the labor market. *Bulletin of LN Gumilyov Eurasian National University. Pedagogy. Psychology. Sociology series.*, 142(1):178–193.

Leili Babashahi, Carlos Eduardo Barbosa, Yuri Lima, Alan Lyra, Herbert Salazar, Matheus Argôlo, Marcos Antonio de Almeida, and Jano Moreira de Souza. 2024. Ai in the workplace: A systematic review of skill transformation in the industry. *Administrative Sciences*, 14(6):127.

Margaret Basson, Tanya Du Plessis, and Roelien Brink. 2023. Visual representation of the mismatch between industry skills demand and higher education skills supply. *International Journal of Work-Integrated Learning*, 24(1):117.

Helena Belchior-Rocha, Inês Casquilho-Martins, and Eduardo Simões. 2022. Transversal competencies for employability: from higher education to the labour market. *Education Sciences*, 12(4):255.

Cedefop. 2023. Setting europe on course for a human digital transition. Accessed June 4, 2025.

Nina Bonderup Dohn, Yasmin Kafai, Anders Mørch, and Marco Ragni. 2022. Survey: Artificial intelligence, computational thinking and learning. *KI-Künstliche Intelligenz*, 36(1):5–16.

A Gelman and J Hill. 2021. Data analysis using regression and multilevel/hierarchical models. 23rd printing. *Analytical Methods for Social Research*.

Marina Johnson, Rashmi Jain, Peggy Brennan-Tonetta, Ethne Swartz, Deborah Silver, Jessica Paolini, Stanislav Mamonov, and Chelsey Hill. 2021. Impact of big data and artificial intelligence on industry: developing a workforce roadmap for a data driven economy. *Global Journal of Flexible Systems Management*, 22(3):197–217.

Douglas C Montgomery, Elizabeth A Peck, and G Geoffrey Vining. 2021. *Introduction to linear regression analysis*. John Wiley & Sons.

Muthmainnah, Prodhan Mahbub Ibna Seraj, and Ibrahim Oteir. 2022. Playing with ai to investigate human-computer interaction technology and improving critical thinking skills to pursue 21st century age. *Education Research International*, 2022(1):6468995.

Davy Tsz Kit Ng, Jac Ka Lok Leung, Samuel Kai Wah Chu, and Maggie Shen Qiao. 2021. Conceptualizing ai literacy: An exploratory review. *Computers and Education: Artificial Intelligence*, 2:100041.

Robert Pater, Herman Cherniaiev, and Marcin Kozak. 2022. A dream job? skill demand and skill mismatch in ict. *Journal of Education and Work*, 35(6-7):641–665.

Michaela Poláková, Juliet Horváthová Suleimanová, Peter Madzík, Lukáš Copuš, Ivana Molnárová, and Jana Polednová. 2023. Soft skills and their importance in the labour market under the conditions of industry 5.0. *Heliyon*, 9(8).

Gursahildeep Singh Sidhu, Md Abu Sayem, Nazifa Taslima, Ahmed Selim Anwar, Fariba Chowdhury, and Manataka Rowshon. 2024. Ai and workforce development: A comparative analysis of skill gaps and training needs in emerging economies. *International journal of business and management sciences*, 4(08):12–28.

Richa Singh Dubey, Justin Paul, and Vijayshri Tewari. 2022. The soft skills gap: a bottleneck in the talent supply in emerging economies. *The International Journal of Human Resource Management*, 33(13):2630–2661.

Yoshija Walter. 2024. Embracing the future of artificial intelligence in the classroom: the relevance of ai literacy, prompt engineering, and critical thinking in modern education. *International Journal of Educational Technology in Higher Education*, 21(1):15.

Jakob Mørup Wang and Zhiru Sun. 2025. LLM-supervised multilingual skill extraction and classification from job ads. In *LNCS, volume 15837*, pages 94–104. Springer Nature.

Xiaojing Weng, Huiyan Ye, Yun Dai, and Oi-lam Ng. 2024. Integrating artificial intelligence and computational thinking in educational contexts: A systematic review of instructional design and student learning outcomes. *Journal of Educational Computing Research*, 62(6):1640–1670.

World Economic Forum. 2025. The future of jobs report 2025: Digest. Accessed June 4, 2025.

# A Appendix



Figure 1: Overview of Our Pipeline Utilizing the Sentence Encoder Twice



Figure 2: Job Volume vs. Skill Demand Over Time

Figure 3: Average Job Demands Across Faculty by Skill Type Over Time



Figure 4: Skill Type Trend by Faculty Over Time

Figure 5: Growing and Declining Skills Across Faculty Over Time



Figure 6: Growing and Declining Skills in Humanities Over Time

# Improving French Synthetic Speech Quality via SSML Prosody Control

**Nassima Ould Ouali**[1]**, Awais Hussain Sani**[2]**, Ruben Bueno**[1†]**,**
**Jonah Dauvet**[1,3†]**, Tim Luka Horstmann**[1,2†]**, Eric Moulines**[1]

[1]École Polytechnique, France,  [2]Hi! PARIS Research Center, France,
[3]McGill University, Canada

{nassima.ould-ouali, ruben.bueno, eric.moulines}@polytechnique.edu

{awais.sani, tim.horstmann}@ip-paris.fr, jonah.dauvet@mail.mcgill.ca

## Abstract

Despite recent advances, synthetic voices often lack expressiveness due to limited prosody control in commercial text-to-speech (TTS) systems. We introduce the first end-to-end pipeline that inserts Speech Synthesis Markup Language (SSML) tags into French text to control pitch, speaking rate, volume and pause duration. We employ a cascaded architecture with two QLoRA-fine-tuned Qwen 2.5-7B models: one predicts phrase-break positions and the other performs regression on prosodic targets, generating commercial TTS-compatible SSML markup. Evaluated on a 14-hour French podcast corpus, our method achieves 99.2% $F_1$ for break placement and reduces mean absolute error on pitch, rate, and volume by 25–40% compared with prompting-only large language models (LLMs) and a BiLSTM baseline. In perceptual evaluation involving 18 participants across over 9 hours of synthesized audio, SSML-enhanced speech generated by our pipeline significantly improves naturalness, with the mean opinion score increasing from 3.20 to 3.87 ($p < 0.005$). Additionally, 15 of 18 listeners preferred our enhanced synthesis. These results demonstrate substantial progress in bridging the expressiveness gap between synthetic and natural French speech. Our code is publicly available at https://github.com/hi-paris/Prosody-Control-French-TTS.

## 1 Introduction

Recent Text-to-Speech (TTS) advances have improved speech intelligibility; yet, achieving natural and expressive prosody remains challenging. Commercial TTS solutions prioritize clarity over prosodic variation, resulting in a monotonous speech output. This limitation particularly affects French due to its complex prosodic features.

**Speech Synthesis Markup Language (SSML)** provides a standardized way to control prosodic features such as pitch, speaking rate, and volume. Unlike neural models, SSML allows post-hoc adjustments and is compatible with commercial TTS engines. Yet, automating SSML generation is difficult: manual markup does not scale, and current LLM-based methods often produce incomplete tags, invalid syntax, or imprecise prosodic control.

We propose an automated SSML pipeline for French, combining structured prosody extraction with a novel cascaded LLM approach for simultaneous tag prediction and prosodic parameter regression. Key contributions include:

- **End-to-end SSML annotation pipeline** that aligns speech to text, segments input into prosodic syntagms, and extracts prosodic coefficients normalized relative to a commercial TTS baseline.

- **Rigorous benchmarking** comparing state-of-the-art (SOTA) approaches (fine-tuned BERT, BiLSTM) with contemporary LLMs across varied prompting strategies and metrics.

- **Cascaded LLM architecture** using two fine-tuned Qwen 2.5-7B models: one for SSML structure/boundaries and another for prosodic prediction, ensuring valid markup and accurate parameter control.

## 2 Related Work

Enhancing neural TTS prosody through automatic markup is an active research domain categorized into: (i) *learning paradigm* (supervised vs. unsupervised approaches) and (ii) *prosodic objective* (prominence, phrasing, style).

**Supervised Prosody Learning**
**Word-level prominence modeling emphasizes salient words using prosodic cues** like pitch and duration. Stephenson et al. (2022) fine-tune BERT (Devlin et al., 2019) to predict three-level

---

†Equal contribution; authors listed in alphabetical order.

prominence tags from wavelet-based labels, achieving $F_1 = 0.588$ and enabling controllable synthesis in FastSpeech 2. Similarly, Zhong et al. (2023) integrate emphasis features into FastSpeech 2, improving expressiveness (+0.49 Mean Opinion Score (MOS)) and naturalness (+0.67 MOS).

**Prosodic emphasis prediction controls automated stress placement patterns.** Shechtman et al. (2021) employ a hybrid model with acoustic and syntactic features, and Seshadri et al. (2021) propose a hierarchical latent model. Liu et al. (2024) combine graph-based contextual encoding with FastSpeech 2 for enhanced rendering. More recently, Chen et al. (2025) present DrawSpeech, a user-sketched prosodic contour control.

**Phrasing segments speech into natural prosodic units with appropriate pauses.** Transformer-based models now outperform recurrent neural networks (RNNs) for break prediction: Futamata et al. (2021) integrate BERT embeddings with linguistic features, improving phrase break prediction ($F_1$ +3.2 points, MOS = 4.39). Vadapalli (2025) show that fine-tuned BERT outperforms RNN baselines, reaching $F_1 = 0.92$ and achieving 58.5% listener preference for BERT-guided punctuation in narrative TTS.

**LLMs enable automated emotional and stylistic annotations at scale**. Yoon et al. (2022) prompt GPT-3 to assign sentence-level emotion labels that guide expressive TTS, achieving MOS 3.92 (naturalness) and 3.94 (expressiveness), matching human-annotated systems. Complementarily, Burkhardt et al. (2023) show that even simple, rule-based SSML adaptations can shape emotional perception, with Unweighted Average Recall scores of 0.76 for arousal and 0.43 for valence.

**Narrative prosody modeling adjusts pitch, speaking rate, and volume to enhance expressive storytelling.** Pethe et al. (2025) use MPNet embeddings and BiLSTMs to predict phrase-level prosody from text. Their SSML-integrated predictions improved alignment with human narration in 22–23 out of 24 audiobooks, yielding +50% listener preference over commercial baselines.

**Unsupervised Prosody Learning**
**Discrete prosody representations eliminate dependency on manual annotations** by learning prosodic patterns directly from speech data. Korotkova et al. (2024) utilize a vector-quantized variational autoencoder with Wav2Vec2 and RoBERTa encodings, deriving ten interpretable prosodic tags

that enhance TTS expressiveness across multiple languages, confirmed by MOS tests ($p < 0.001$). In contrast, Karlapati et al. (2021) learn continuous 64-dimensional prosody embeddings: a VAE encodes mel-spectrograms, and a RoBERTa + syntax-GNN regresses these from text. At inference, the 64-dimensional prosodic code conditions a Tacotron2 decoder, yielding 13.2% comparative MOS gain ($3.30 \rightarrow 3.74$, $p < 0.005$) on LJSpeech with $F_0$ correlation of $r = 0.68$. Discrete tags offer interpretability; continuous embeddings better capture fine-grained intonation. Both improve TTS expressiveness without hand-crafted annotations.

**Limitations of Prior Work:** However, existing research exhibits critical gaps. Current methods lack a comprehensive end-to-end framework for converting raw speech into standardized SSML-compliant prosodic markup. Most rely on partial manual annotations, address isolated prosodic control aspects, or produce markup incompatible with commercial TTS systems. Furthermore, the majority of existing work focuses on English, leaving other complex languages like French under-explored. Additionally, current LLM-based approaches suffer from systematic limitations, undergenerating necessary tags, producing syntactically invalid SSML structures, and lacking precise control over numerical prosodic parameters, which prevents deployment in practical TTS systems.

We address these limitations with two main contributions: (i) we introduce the first reproducible, comprehensive French pipeline that automatically extracts fine-grained prosodic annotations and converts them into standards-compliant SSML, and (ii) we develop a novel cascaded LLM architecture that generates syntactically correct prosodic tags with precise numerical control at inference time, resulting in substantially enhanced naturalness and expressiveness in synthetic speech.

## 3 Dataset Creation and SSML Annotation Pipeline

We construct a comprehensive dataset annotated with prosodic features from French speech. Our methodology involves aligning spoken audio with transcripts, extracting four key prosodic features — **pitch**, **volume**, **speaking rate**, and **break duration** — and converting them into standardized SSML for enhanced synthetic speech generation. Figure 1 presents our preprocessing pipeline. Further dataset statistics are provided in Appendix A.

Figure 1: Overview of the SSML annotation pipeline. Natural speech is aligned, segmented, and compared to a synthetic baseline to extract prosodic features for SSML markup. Green elements indicate later model training data.

**Audio Collection and Preprocessing:** We process 14 hours of diverse French audio content sourced from *ETX Majelan*[1], a high-quality podcast platform with interviews and discussions. Our dataset includes speech from 14 distinct speakers (42% female). The original recordings contain background music, jingles, and sound effects, complicating prosodic analysis. Hence, we isolate clean speech using *Demucs* (Défossez, 2021), a SOTA audio source separation model, down-sample to 16 kHz, and peak-normalize the audio. Using *pydub*[2], we segment the cleaned audio via silence detection with a $-35$ dBFS threshold and 300 ms gaps. The resulting audio segments serve as our fundamental processing units for subsequent prosodic analysis.

**Text-Audio Alignment:** Accurate alignment between audio and transcribed text is crucial for prosody extraction, but particularly challenging in French due to phonetic phenomena such as liaisons, elisions, and prosodic contractions. To address this, we employ the Whisper Timestamped package[3] with the Whisper (Radford et al., 2022) medium model and Auditok Voice Activity Detection (VAD)[4], which filters out silent segments that would otherwise distort prosodic measurements.

To evaluate this setup, we benchmarked it against larger Whisper models, Montreal Forced Aligner (MFA) by McAuliffe et al. (2017), and NVIDIA NeMo (Kuchaiev et al., 2019). Benchmarking used our dataset and FLEURS benchmark (Conneau et al., 2022) as a state-of-the-art reference.. While larger models yielded marginal gains, they introduced instability such as significant hallucinations during silence – a known issue in Whisper (Barański et al., 2025) – as well as higher computational cost. Our chosen configuration

achieved a 5.95% WER using Whisper-medium with an average Alignment Recall Rate (ARR) of 96.3% over 15-second windows against the manual TextGrid annotations created with Praat (Boersma and Van Heuven, 2001) (see Table 1). This yielded an optimal accuracy-efficiency trade-off.

**Baseline Voice for Prosodic Comparison:** For prosodic reference, we synthesize each transcript using `Microsoft Azure Neural TTS` with the French voice *Henri* (Microsoft Azure, 2024). Henri was selected for its clarity, broad phonetic coverage, and consistent yet neutral prosodic characteristics, making it optimal for computing relative prosodic adjustments. The resulting synthetic speech provides a stable baseline against which natural prosodic features are measured and compared, as detailed in subsequent sections.

**Syntagm Segmentation:** Each segment undergoes further subdivision into *syntagms*: prosodic units with natural pause boundaries. Following Roll et al. (2023), we detect them through acoustic pauses and punctuation. We first derive a word/pause sequence from the TextGrid, where pauses following function words are discarded with a POS filter to remove Whisper artifacts. Next, any silence that follows ., *?*, or *!* is clamped to at least 500 ms, and a 500 ms pause is injected whenever Whisper failed to signal the end of a sentence. The resulting timestamped syntagms provide stable, linguistically meaningful units for prosodic analysis.

**Prosodic Feature Extraction and SSML Tag Construction:** Each syntagm is annotated with four prosodic features: median **pitch** (fundamental frequency $f_0$), segment-level **volume** (Loudness Units Full Scale (LUFS)), **speaking rate** (words per second), and inter-syntagmatic **break duration**. All features are computed for both natural and synthetic baseline voices to derive relative delta values for SSML encoding. To account for intra-

---

[1] https://etxmajelan.com/
[2] https://github.com/jiaaro/pydub
[3] https://github.com/linto-ai/whisper-timestamped
[4] https://github.com/amsehili/auditok

Table 1: Metric evaluation of four whisper models, MFA, and NeMo on our dataset (Section 3) and FLEURS.

| Metric | Whisper Model Variants | | | | Alignment Models | |
|---|---|---|---|---|---|---|
| | Medium | Large v2 | Large v3 | Turbo Large v3 | MFA ‡ | NeMo Large ‡ |
| Parameters | 769 M | 1550 M | 1540 M | 809 M | – | 120 M |
| WER† | 5.95% / 10.70% | 4.60% / 6.27% | 3.92% / 5.80% | 3.52% / 5.71% | – | – |
| WER† +VAD | 5.68% / 8.72% | 5.07% / 6.16% | 3.86% / 5.65% | 6.16% / 5.83% | – | – |
| ARR* | 96.3% | 97.1% | 96.2% | 97.8% | 99.7% | 50.7% |
| Start MAE* (ms) | 264 | 191 | 207 | 152 | 115 | 4529 |
| Duration MAE* (ms) | 91 | 77 | 102 | 76 | 95 | 218 |

† WER computed with the HuggingFace `evaluate` library.
∗ ARR and MAE are computed on 15 second segments, against the gold manually annotated Text Grids of 1 hour of speech from our dataset.
‡ MFA and NeMo alignments use gold transcripts, thus rendering the WER 0 by default.

and inter-speaker variability, we normalize each syntagm's pitch, volume, and rate relative to a baseline computed as the median over a sliding window of $w = 10$ audio segments (or, when $w$ covers all segments, the global median). The computation of each feature is detailed as follows:

**Pitch** median fundamental frequency $f_0^{(i)}$ is converted to a semitone offset $s_i = 12 \log_2\big(f_0^{(i)}/\bar{f}_0\big)$, clipped to $\big[-0.7P, \ P\big]$ to allow larger upward than downward shifts, and re-scaled to percentage pitch change $p_i = \big(2^{s_i/12} - 1\big) \times 100$.

Using LUFS for **volume**, the baseline–synthetic difference $\Delta L_i = \bar{L} - L_{\text{syn}}^{(i)}$ is mapped to a gain $v_i = \big(10^{\Delta L_i/20} - 1\big) \times 100$, then clipped to $\pm V$ (we use $V = 10\%$).

**Speaking rate** is estimated as *words per second*. Let $n_i$ be the word count and $d_{\text{nat}}$, $d_{\text{syn}}$ the net speaking durations (pauses removed). The rate delta is $r_i = \frac{n_i/d_{\text{nat}} - n_i/d_{\text{syn}}}{n_i/d_{\text{syn}}} \times 100$. Slow-downs are amplified for long syntagms ($> 1$ s) while speedups are reduced, and the final value is clamped to $\pm R$ with a tighter $+0.5R$ ceiling for accelerations.

To improve prosodic smoothness, we apply exponential smoothing to pitch and rate with $\alpha = 0.2$:

$$\tilde{x}_0 = x_0, \quad \tilde{x}_i = \alpha x_i + (1 - \alpha)\tilde{x}_{i-1}.$$

Sudden jumps are clamped to $\Delta = 8\%$ per syntagm. Volume is not smoothed.

**Break durations** are taken from the intersyntagm silence gaps and inserted as raw durations (e.g., `<break time="200ms">`). The final SSML markup is assembled by inserting appropriate `<prosody>` and `<break>` tags into the text.[5]

---

[5]At inference, we found that wrapping each `<prosody>` tag with `<mstts:silence type="leading-exact/trailing-exact" value="0"/>` improves output by suppressing unwanted Azure TTS pauses.

## 4 Methodology

We test whether text alone encodes sufficient cues for prosody by training two baselines: (i) a BERT-base model fine-tuned for token-level pause prediction (Vadapalli, 2025), and (ii) a BiLSTM (Pethe et al., 2025) which predicts SSML tags with pitch, speaking rate and volume adjustments.

### 4.1 Fine-tuning BERT for Pause Prediction

Following Vadapalli (2025), we fine-tune an uncased BERT-base model for token-level pause prediction. A binary classification head determines whether each sub-word is followed by a break tag. We adopt the same hyperparameters as the original work on our dataset: batch size 64, learning rate $10^{-5}$, and gradient clipping at 10. For evaluation, we report both $F_1$ score and perplexity. While $F_1$ is used in Vadapalli (2025), perplexity is introduced here as an additional metric to enable broader comparisons in later sections.

We additionally introduce bootstrapping, a technique not used in the original paper, to evaluate the small model's variance in performance. We bootstrap on 10 distinct sets with the same configuration as the original training set, which allows us to obtain a distribution of performance scores for robust estimation of the uncertainty of performance. Given the reduced size of the dataset, we expect overall performance to degrade slightly. Hence, we focus on stability, measured via standard deviation.

### 4.2 BiLSTM-Based Sequence Modeling

We implement a BiLSTM baseline following Pethe et al. (2025), explicitly modeling prosody prediction as a sequence regression task to predict three SSML parameters: pitch, volume, and speaking rate. This approach leverages local context through sequential processing of prosodic units.

Each syntagm receives encoding into a 768-dimensional representation using the pretrained sentence encoder `all-mpnet-base-v2`[6]. We construct overlapping input sequences of varying lengths $L \in \{1, 2, 3, 4\}$, extending beyond the original study's sequence lengths of 2 and 3 to assess optimal context window size. targeting z-scored prosody vectors (pitch, volume, rate) of central segments, normalized on training statistics.

The architecture includes LayerNorm preprocessing, bidirectional LSTM (40 units per direction), dense layer (20 units, $\tanh$ activation), and linear projection for predicting the 3-dimensional prosody vector. Training uses MSE loss between predicted and target z-scored vectors. We additionally compute raw RMSE and MAE metrics for interpretability and literature comparison.

### 4.3 Zero-shot and Few-shot Evaluation

To assess SOTA LLMs for SSML markup generation, we benchmarked various open-source models in both **zero-shot** and **few-shot** settings. We evaluated Mistral (7B), Qwen 2.5 (7B), Llama 3 (8B), Granite 3.3 (8B), Qwen 3 (8B), DeepSeek-R1 (32B), and Qwen 3 (32B) via the Ollama framework[7]. Models were prompted at the segment level ($\approx$ eight tags per segment on average) with French text, and tasked with generating fully annotated SSML for 100 randomly chosen segments. Few-shot prompts included 10 reference examples.

### 4.4 Cascaded Fine-tuning Approach

As we show in Section 5.3, LLM-based approaches under-generate `<break>` and `<prosody>` tags, resulting in SSML that is structurally incomplete and limited in expressive control. To address this, we introduce a cascaded strategy that separates structural and numerical prediction. The first model, *QwenA*, predicts where prosodic boundaries occur; the second, *QwenB*, supplies the corresponding numerical attributes.

**QwenA (Stage 1): Break Prediction**

We fine-tune a Qwen 2.5-7B model (QLoRA with 4-bit quantization, rank 8, $\alpha = 16$) to insert `<break>` tags at linguistically appropriate junctures. QwenA processes up to 200-word French paragraphs (within a 1024-token limit), retaining punctuation, quotations, and parenthetical clauses so

that the model must reason over long-range dependencies rather than relying on sentence-level heuristics. These features reflect real-world TTS applications, where systems rarely receive inputs entirely devoid of punctuation or other natural orthographic cues. Furthermore, this approach aligns with the baseline methodology used in Vadapalli (2025). A deterministic post-processor then converts each `<break>` into an empty `<prosody>` element, yielding a syntactically valid SSML skeleton to pass into the next stage.

**QwenB (Stage 2): Prosodic Regression**

QwenB builds on the skeleton emitted by QwenA and replaces each empty `<prosody>` placeholder with fully specified numeric attributes (pitch, rate, volume, and break duration). Starting again from Qwen 2.5-7B, we inject a second QLoRA adapter with 4-bit quantization (rank 8, $\alpha = 16$) into the value and feed-forward projections so that only those low-rank updates are trainable. **Loss is computed on the numeric tokens**, so categorical text incurs zero penalty and the adapter's capacity is **devoted entirely to modeling prosodic distributions**. Targets are standardized to unit variance during optimization and rescaled at inference, a choice that stabilizes gradients and accelerates convergence.

## 5 Results and Analysis

### 5.1 Perceptual Evaluation (AB Test)

To assess our SSML annotation pipeline's effectiveness in enhancing synthetic speech (Section 3), we conducted AB testing with 18 participants. Each participant evaluated 30 one-minute audio pairs, where the baseline was the raw, unaltered voice of Microsoft Azure Neural TTS (Henri), without any prosody modifications, compared to the prosody-enhanced version. These pairs were presented randomly, with 60 segments evaluated per participant.

The SSML-enhanced audio achieved a MOS of 3.87 (5-point scale), outperforming the baseline (3.20) and yielding a **20% improvement** in perceived quality. Additionally, 15 of 18 participants preferred the enhanced version in over half of the cases, with 7 preferring it in more than 75% of comparisons. These results support the effectiveness of our SSML-based prosody enhancement for improving synthetic speech quality.

---

[6] https://huggingface.co/sentence-transformers/all-mpnet-base-v2

[7] https://ollama.com/

Figure 2: Cascaded LLM approach for automated text-to-SSML generation: QwenA predicts tag placement, QwenB injects prosodic values. This disentangled design enables accurate and efficient prosody control for synthetic speech.

## 5.2 Baseline Model Performance

**BERT Break Prediction Results:** We evaluated the performance of the fine-tuned BERT model from Vadapalli (2025) on $F_1$ (%) and perplexity (best = 1), and attained results very close to the original paper, which reports a 92.10% $F_1$ score for break prediction. Our model achieves a $F_1$ score of 92.06%, along with a perplexity of 1.123 (not reported in Vadapalli (2025) but useful for further evaluation). Our stability assessment on 10 bootstrapped datasets yielded an average $F_1$ of 47.52% $\pm$ 4.65% (Confidence Interval (CI) = 9.8%) and perplexity of $1.274 \pm 0.005$ (CI = 0.4%), indicating high stability in token prediction but moderate variability in classification performance. We present the results from the original training data in Table 4.

**BiLSTM Prosody Prediction Results:** We evaluated our BiLSTM model following Pethe et al. (2025). Table 2 presents the MSE values for normalized z-score prosody features. Our approach achieves SOTA results comparable to those reported in the original paper. For a more comprehensive analysis, we also report the raw RMSE and MAE (%) for each prosodic parameter.

Unlike Pethe et al. (2025), our analysis revealed that a sequence window length ($L = 2$) yielded superior performance across prosodic attributes. Specifically, $L = 2$ demonstrated lower error rates for two of the three prosodic attributes, while pitch prediction achieved optimal performance with $L = 1$. Notably, the MAE for volume with $L = 2$ was more than 0.04 percentage points lower than all other tested lengths, and 0.04–0.09 percentage points superior to alternative sequence configurations. Our best results align with those of Pethe et al. (2025): z-scored MSE of 0.8734 for

Table 2: BiLSTM-based prosodic attribute prediction across sequence window lengths ($L$). Best overall performance is achieved at $L = 2$, with lowest MAE for volume and rate, and near-best scores for pitch.

| $L$ | Metric Type | Pitch | Volume | Rate |
|---|---|---|---|---|
| 1 | Z-score MSE ($\downarrow$) | 0.8752 | 0.9141 | 0.7733 |
| | % RMSE ($\downarrow$) | 2.0659 | 7.8597 | 1.2771 |
| | % MAE ($\downarrow$) | 1.6709 | 6.4768 | 0.8878 |
| 2 | Z-score MSE ($\downarrow$) | 0.8983 | 0.8949 | 0.7572 |
| | % RMSE ($\downarrow$) | 2.0930 | 7.7767 | 1.2637 |
| | % MAE ($\downarrow$) | 1.6883 | 6.0405 | 0.8462 |
| 3 | Z-score MSE ($\downarrow$) | 0.9936 | 0.9917 | 0.8593 |
| | % RMSE ($\downarrow$) | 2.2012 | 8.1864 | 1.3462 |
| | % MAE ($\downarrow$) | 1.7732 | 6.5100 | 0.9257 |
| 4 | Z-score MSE ($\downarrow$) | 0.9950 | 0.9992 | 0.8263 |
| | % RMSE ($\downarrow$) | 2.2028 | 8.2172 | 1.3201 |
| | % MAE ($\downarrow$) | 1.7568 | 6.5990 | 0.9312 |

pitch, 0.7631 for volume, and 1.0610 for speaking rate. We attribute minor performance differences to dataset variations and establish the $L = 2$ model as our primary baseline for subsequent comparisons with our proposed cascaded architecture.

## 5.3 Zero-shot and Few-shot Prompting Evaluation

We first focused on evaluating *break tag prediction*, a proxy for assessing structural correctness and syntagm segmentation. Figure 3 shows the average number of predicted <break> and <prosody> tags per segment compared to gold annotations. All models consistently under-generate tags, indicating systematic issues maintaining SSML structure. Few-shot prompting led to unexpected patterns: fewer predicted break tags but increased <prosody> tags, suggesting attention shifts or

Table 3: SSML generation performance across models and prompting strategies, evaluated by cosine similarity of predicted vs. gold SSML embeddings, and MAE/RMSE for pitch, volume, rate, and break durations. Qwen2.5 (7B) offers the best trade-off between accuracy and efficiency.

| Model | SSML Sim. ↑ | Pitch (%) MAE/RMSE↓ | Volume (%) MAE/RMSE↓ | Rate (%) MAE/RMSE↓ | Break Time (ms) MAE/RMSE↓ |
|---|---|---|---|---|---|
| Qwen3 (32B) (ZS) | 0.91 | 1.42/1.83 | 7.65/8.48 | 1.52/2.00 | 170.23/232.41 |
| Qwen2.5 (7B) (ZS) | 0.90 | 2.07/2.43 | 7.23/8.05 | 1.54/1.93 | 361.88/393.04 |
| Qwen3 (32B) (FS) | 0.90 | 1.08/1.41 | 5.80/7.33 | 0.97/1.31 | 159.58/215.50 |
| Qwen3 (8B) (FS) | 0.90 | 1.77/2.83 | 6.96/16.85 | 1.23/1.69 | 147.24/242.98 |
| Qwen2.5 (7B) (FS) | 0.89 | 1.26/1.50 | 4.32/6.77 | 1.01/1.24 | 118.85/179.68 |
| Mistral (7B) (ZS) | 0.88 | 1.85/2.25 | 24.19/43.96 | 18.30/41.24 | 207.28/258.76 |
| Mistral (7B) (FS) | 0.87 | 1.75/2.16 | 5.38/8.33 | 1.14/1.42 | 205.03/384.17 |
| Granite3.3 (8B) (FS) | 0.87 | 1.45/1.86 | 4.95/7.12 | 0.95/1.30 | 196.93/265.07 |
| Llama3 (8B) (ZS) | 0.84 | 1.44/1.82 | 7.30/8.08 | 2.26/10.17 | 285.17/318.19 |
| Qwen3 (8B) (ZS) | 0.82 | 1.99/2.70 | 7.43/8.41 | 1.69/2.06 | 274.27/334.20 |
| Deepseek-R1 (32B) (ZS) | 0.81 | 1.64/2.11 | 15.50/30.41 | 18.79/41.14 | 274.66/320.62 |
| Granite3.3 (8B) (ZS) | 0.76 | 3.70/4.55 | 13.86/29.11 | 33.25/55.85 | 320.77/413.91 |
| Deepseek-R1 (32B) (FS) | 0.76 | 1.43/2.04 | 7.12/8.23 | 3.69/12.94 | 244.85/302.87 |
| Llama3 (8B) (FS) | 0.34 | 1.26/1.62 | 7.24/8.23 | 1.53/1.88 | 416.13/445.99 |

↑: higher is better, ↓: lower is better. ZS: Zero-Shot, FS: Few-Shot.



(a) Break tag usage comparison (DS = DeepSeek)



(b) Prosody tag usage comparison (DS = DeepSeek)

Figure 3: Structural comparison of SSML tag predictions across models. All models under-generate both break and prosody tags relative to the gold standard.

stylistic overfitting to prompt exemplars. Notably, Llama 3 and DeepSeek-R1 (32B) show large discrepancies between zero- and few-shot modes, with Llama 3's prosody tagging almost collapsing in the few-shot case.

Beyond structural accuracy, we evaluated numerical performance through *cosine similarity* between predicted and reference SSML structures, embedded using the *all-MiniLM-L6-v2* [8] model, RMSE, and MAE for break durations and prosodic coeffi-

cients, averaged per segment. Table 3 summarizes the results. Qwen 2.5-7B achieves the best overall balance: in the few-shot setting, it delivers the lowest MAE for break (118.85 ms) and volume (4.32%), and second-highest structural similarity in zero-shot (0.9). Qwen 3 (32B) slightly surpasses it on similarity (0.908), but at a cost of 4.5 times higher memory usage and slower inference, making it less suitable for fine-tuning and deployment.

Our findings suggest that while few-shot prompting can improve prosody tag usage and numerical accuracy, model behavior is highly architecture-dependent. Furthermore, the consistent underproduction of tags across models highlights the need for more robust SSML-structure awareness.

## 5.4 Cascaded LLM Evaluation

Our evaluation of the cascaded QwenA and QwenB models demonstrates substantial performance improvements over existing SOTA approaches, as detailed in Table 4:

Table 4: Break tag prediction: F1 and perplexity for our cascaded model (QwenA) vs. fine-tuned BERT. QwenA achieves near-perfect accuracy and fluency.

| Model | F1 (%) | Perplexity ($\rightarrow$ 1) |
|---|---|---|
| Cascade (QwenA) | 99.24 | 1.00 |
| Finetuned BERT | 92.06 | 1.12 |

For QwenA, which utilizes next-token prediction on a linearized SSML target, the model achieved a test **perplexity of 1.001 and a tag-level $F_1$ score**

Table 5: RMSE (↓) and MAE (↓) for our cascaded model vs. benchmarks. It achieves the lowest error scores across nearly all prosody attributes.

| Model | Metric | Pitch | Volume | Rate | Break Time |
|---|---|---|---|---|---|
| Cascade (Ours) | RMSE | 1.22 | 1.67 | 1.50 | 166.51 |
| | MAE | 0.97 | 1.09 | 1.10 | 132.89 |
| BiLSTM[†] ($L = 2$) | RMSE | 2.09 | 7.77 | 1.26 | – |
| | MAE | 1.68 | 6.04 | 0.84 | – |
| SOTA Few-Shot* | RMSE | 1.41 | 7.33 | 1.31 | 215.50 |
| | MAE | 1.08 | 5.80 | 0.97 | 159.58 |
| SOTA Zero-Shot* | RMSE | 1.83 | 8.48 | 2.00 | 232.41 |
| | MAE | 1.42 | 7.65 | 1.52 | 170.23 |

†: Results based on Pethe et al. (2025); see Section 4.2
∗: Qwen-3 (32B) selected via cosine similarity (Tab. 3)
Units: Pitch, Volume, Rate (%); Break Time (ms)

**of 99.24%**, surpassing the fine-tuned BERT's baseline of 1.123 perplexity and 92.06% $F_1$ score. Moreover, this approach also outperforms the LLM tag prediction benchmarks, which consistently under-generate break and prosody tags, as illustrated in Figures 3a and 3b. This near-perfect tag insertion accuracy validates the improved performance of our cascaded approach compared to available models for SSML tag prediction.

QwenB demonstrates significant advancements in prosody parameter prediction, achieving an **MAE of 0.97% for pitch, 1.09% for volume, 1.10% for rate, and 132.9ms for break timing** (Table 5). Furthermore, this strong performance is achieved while maintaining an efficient end-to-end latency of approximately 190 ms for a 150-word paragraph. This demonstrates the model's enhanced SSML parameter prediction and its ability to process larger text segments, outperforming baseline approaches. This performance also suggests that evaluations of pipeline audio (Section 5.1) are highly generalizable to the cascaded model's audio quality due to their close similarity.

### 5.5 Summary and Analysis of Results

Table 5 provides a comparative overview of objective performance across all evaluated approaches, revealing three key observations:

1. **Cascaded QwenA + QwenB sets new SOTA performance.** The system achieves single-digit MAE for all prosodic coefficients and reduces break-timing error by 25% vs. the best few-shot LLM baseline.

2. **BiLSTM architectures remain competitive for speaking rate prediction.** Though outperformed elsewhere, its 0.84% MAE on rate shows lightweight sequential models still capture localized prosodic patterns effectively.

3. **Prompt-only LLMs systematically under-generate tags.** Both zero- and few-shot settings underperform supervised baselines on break timing prediction (MAE > 150 ms) and structural metrics (Figure 3), reinforcing the necessity for explicit structural supervision in SSML generation tasks.

These findings confirm that disentangling structural prediction (QwenA) from numerical regression (QwenB) yields optimal performance across both dimensions: syntactically valid SSML markup with fine-grained prosodic control, while preserving inference efficiency suitable for real-time TTS applications. The subjective evaluation results in Section 5.1 corroborate these objective improvements, demonstrating that enhanced technical performance translates into substantial perceptual gains—a 20% MOS improvement and consistent listener preference for enhanced synthesis.

## 6 Conclusion and Future Directions

Using a fine-tuned cascaded Qwen 2.5-7B architecture, we separate structural tag insertion from prosodic parameter prediction, achieving near-perfect break placement (99.2% $F_1$, perplexity 1.001) and reducing prosodic MAE below 1.1 points – representing 25–40% better than prompting-only LLMs and BiLSTM baselines.

Perceptual evaluation shows that SSML from our pipeline increases MOS from 3.20 to 3.87, with consistent listener preference. This marks a significant step toward closing the expressiveness gap between synthetic and natural French speech while preserving compatibility with commercial TTS.

Future research includes unifying our cascaded approach into a single end-to-end model for joint prosodic prediction, incorporating multimodal audio embeddings to capture subtle speech characteristics beyond text-derived features, and extending this methodology to additional languages to assess cross-linguistic generalizability and robustness.

## 7 Limitations

While our proposed system shows significant improvements, several limitations warrant discussion. Our experiments focus exclusively on French using a proprietary 14-hour corpus. While our pipeline is language-agnostic, performance may vary for languages with different prosodic characteristics. The dataset size remains modest compared to typical TTS training corpora, as English prosody modeling often leverages hundreds of hours of annotated speech, indicating that scaling our French dataset could yield additional performance gains. Additionally, our improvements rely on TTS engines supporting fine-grained SSML tags, meaning legacy or non-compliant systems may not achieve similar gains and may require custom adjustments for engine-specific behaviors.

Our prosodic deltas are computed with respect to a single baseline synthetic voice (Azure fr-FR-HenriNeural) and evaluated with the same voice, which limits out-of-domain generalization. While SSML prosody tags are standardized, their acoustic realization is implementation- and voice-dependent; engines may clamp or substitute values, and different voices can map the same percentage to different F0/rate changes. Consequently, numeric SSML settings may require voice-specific recalibration (e.g., a short script that sweeps pitch/rate/volume and measures resulting semitone, syllables/s, and dB changes) before transfer to other voices or engines.

From a computational perspective, fine-tuning Qwen 2.5-7B requires substantial GPU memory ($\approx 15\ GB$ peak) despite 4-bit quantization, necessitating model compression or distillation for smaller deployments. Conversely, greater computational resources could enable more extensive fine-tuning and potentially improve performance. Our approach also assumes that punctuation and syntactic cues correlate well with natural prosodic boundaries, an assumption that may break down in highly informal or unpunctuated text such as social media transcripts, leading to suboptimal break placement.

## 8 Ethics Statement

Our work uses commercially licensed French podcast audio, ensuring no personal or sensitive data are exposed. We acknowledge potential biases from using a limited speaker set and encourage broader demographic validation. While improved prosody can enhance synthetic voices, it also risks misuse in deceptive audio generation; we therefore recommend watermarking or verification mechanisms. Code and anonymized alignment scripts are publicly shared to promote reproducibility and transparency.

## References

Mateusz Barański, Jan Jasiński, Julitta Bartolewska, Stanisław Kacprzak, Marcin Witkowski, and Konrad Kowalczyk. 2025. Investigation of Whisper ASR Hallucinations Induced by Non-Speech Audio. In *ICASSP 2025 - 2025 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, pages 1–5. ArXiv:2501.11378 [cs].

Paul Boersma and Vincent Van Heuven. 2001. Speak and unspeak with praat. *Glot Int*, 5:341–347.

Felix Burkhardt, Uwe Reichel, Florian Eyben, and Björn Schuller. 2023. Going retro: Astonishingly simple yet effective rule-based prosody modelling for speech synthesis simulating emotion dimensions. *arXiv preprint arXiv:2307.02132*.

Estelle Campione and Jean Véronis. 2002. A large-scale multilingual study of silent pause duration. In *Speech Prosody 2002*, pages 199–202. ISCA.

Weidong Chen, Shan Yang, Guangzhi Li, and Xixin Wu. 2025. Drawspeech: Expressive speech synthesis using prosodic sketches as control conditions. *arXiv preprint arXiv:2501.04256*.

Alexis Conneau, Min Ma, Simran Khanuja, Yu Zhang, Vera Axelrod, Siddharth Dalmia, Jason Riesa, Clara Rivera, and Ankur Bapna. 2022. Fleurs: Few-shot learning evaluation of universal representations of speech.

Alexandre Défossez. 2021. Hybrid spectrogram and waveform source separation. In *Proceedings of the ISMIR 2021 Workshop on Music Source Separation*.

Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. 2019. BERT: Pre-training of Deep Bidirectional Transformers for Language Understanding. ArXiv:1810.04805 [cs].

Kosuke Futamata, Byeongseon Park, Ryuichi Yamamoto, and Kentaro Tachibana. 2021. Phrase break prediction with bidirectional encoder representations in japanese text-to-speech synthesis. *arXiv preprint arXiv:2104.12395*.

Sri Karlapati, Ammar Abbas, Zack Hodari, Alexis Moinet, Arnaud Joly, Penny Karanasou, and Thomas Drugman. 2021. Prosodic representation learning and contextual sampling for neural text-to-speech. In *ICASSP 2021-2021 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, pages 6573–6577. IEEE.

Yuliya Korotkova, Ilya Kalinovskiy, and Tatiana Vakhrusheva. 2024. Word-level text markup for prosody control in speech synthesis. In *Proc. Interspeech 2024*, pages 2280–2284.

Oleksii Kuchaiev, Jason Li, Huyen Nguyen, Oleksii Hrinchuk, Ryan Leary, Boris Ginsburg, Samuel Kriman, Stanislav Beliaev, Vitaly Lavrukhin, Jack Cook, Patrice Castonguay, Mariya Popova, Jocelyn Huang, and Jonathan M. Cohen. 2019. NeMo: a toolkit for building AI applications using Neural Modules. ArXiv:1909.09577 [cs].

Rui Liu, Zhenqi Jia, Jie Yang, Yifan Hu, and Haizhou Li. 2024. Emphasis rendering for conversational text-to-speech with multi-modal multi-scale context modeling. *arXiv preprint arXiv:2410.09524*.

Michael McAuliffe, Michaela Socolof, Sarah Mihuc, Michael Wagner, and Morgan Sonderegger. 2017. Montreal Forced Aligner: Trainable Text-Speech Alignment Using Kaldi. In *Interspeech 2017*, pages 498–502. ISCA.

Microsoft Azure. 2024. Speech synthesis markup language (ssml) documentation.

Naomi Peck and Laura Becker. 2024. Syntactic pausing? Re-examining the associations. *Linguistics Vanguard*, 10(1):223–237. Publisher: De Gruyter Mouton.

Charuta Pethe, Bach Pham, Felix D Childress, Yunting Yin, and Steven Skiena. 2025. Prosody analysis of audiobooks.

Alec Radford, Jong Wook Kim, Tao Xu, Greg Brockman, Christine McLeavey, and Ilya Sutskever. 2022. Robust speech recognition via large-scale weak supervision. *arXiv preprint arXiv:2212.04356*.

Nathan Roll, Calbert Graham, and Simon Todd. 2023. PSST! prosodic speech segmentation with transformers. In *Proceedings of the 27th Conference on Computational Natural Language Learning (CoNLL)*, pages 476–487, Singapore. Association for Computational Linguistics.

Shreyas Seshadri, Tuomo Raitio, Dan Castellani, and Jiangchuan Li. 2021. Emphasis control for parallel neural tts. *arXiv preprint arXiv:2110.03012*.

Slava Shechtman, Raul Fernandez, and David Haws. 2021. Supervised and unsupervised approaches for controlling narrow lexical focus in sequence-to-sequence speech synthesis. In *2021 IEEE Spoken Language Technology Workshop (SLT)*, pages 431–437.

Brooke Stephenson, Laurent Besacier, Laurent Girin, and Thomas Hueber. 2022. BERT, can HE predict contrastive focus? Predicting and controlling prominence in neural TTS using a language model. ArXiv:2207.01718 [cs].

Anandaswarup Vadapalli. 2025. An investigation of phrase break prediction in an end-to-end tts system. *SN Computer Science*, 6(2):1–11.

Hyun-Wook Yoon, Ohsung Kwon, Hoyeon Lee, Ryuichi Yamamoto, Eunwoo Song, Jae-Min Kim, and Min-Jae Hwang. 2022. Language model-based emotion prediction methods for emotional speech synthesis systems. *arXiv preprint arXiv:2206.15067*.

Yi Zhong, Chen Zhang, Xule Liu, Chenxi Sun, Weishan Deng, Haifeng Hu, and Zhongqian Sun. 2023. Ee-tts: Emphatic expressive tts with linguistic information. *arXiv preprint arXiv:2305.12107*.

# A  Dataset Statistics

The dataset constructed through our end-to-end SSML annotation pipeline (Section 3) comprises 14 speakers (42% female), encompassing 122,303 words across 711,603 characters. Our annotation process generated 17,695 `<prosody>` tags and 18,746 `<break>` tags, providing comprehensive prosodic markup for the corpus (Table 6).

Table 6: Corpus statistics for the annotated French speech dataset

| Metric | Value |
|---|---|
| Speakers | 14 |
| Total characters | 711,603 |
| Total words | 122,303 |
| Prosody tags | 17,695 |
| Break tags | 18,746 |

The prosodic parameter distributions reveal linguistically meaningful patterns (Figure 4). Pitch adjustments cluster around +1% with 50% of values within ±2%, reflecting the subtle phrase-final rises characteristic of French declarative intonation. Rate modifications center at -1%, indicating slight deceleration relative to the neutral Azure baseline, consistent with the deliberate pacing typical of podcast narration. Volume adjustments concentrate at -10% with an upper bound at +2%, reflecting our systematic reduction strategy relative to the synthetic baseline to achieve more natural amplitude levels.



Figure 4: Distribution of prosodic parameters in the annotated dataset. **Left:** Pitch, rate, and volume adjustments (percentage) relative to synthetic baseline. **Right:** Break durations (milliseconds) derived from natural inter-phrasal pauses.

Break duration analysis reveals a median pause of approximately 400 ms with an interquartile range of 250–500 ms, aligning with established phonetic studies on French prosodic phrase boundaries (Peck and Becker, 2024; Campione and Véronis, 2002).

# B  Comparative Analysis of Prosodic Features

## B.1  Pitch Characteristics

Figure 5 demonstrates the temporal evolution of fundamental frequency in natural versus synthesized speech. Natural speech exhibits a broader pitch range with complex, fluid intonational patterns reflecting the dynamic modulation inherent in human vocal production. Conversely, synthesized speech operates within a constrained, generally lower fundamental frequency range, displaying more abrupt transitions and reduced prosodic variability.



Figure 5: Temporal pitch contours, comparing natural and synthesized speech across representative utterances

Cross-speaker analysis (Figure 6) reveals substantial inter-speaker pitch variability in natural speech, while synthesized versions cluster within a significantly narrower frequency range. This compression of the pitch space in synthetic speech represents a fundamental limitation in current TTS systems' ability to capture individual vocal characteristics.

## B.2  Volume Dynamics

Amplitude modulation patterns (Figure 7) reveal marked differences between natural and synthetic speech production. Natural speech demonstrates substantial dynamic range with frequent amplitude variations, characteristic of expressive human discourse and reflecting the speaker's communicative intent. Synthesized speech exhibits limited volume

Figure 6: Speaker-wise mean pitch comparison: natural speech (y-axis) versus synthesized speech (x-axis). Each point represents one speaker's average fundamental frequency.

variation, maintaining relatively consistent amplitude levels that contribute to reduced prosodic expressiveness.



Figure 7: Volume variation patterns over time for natural versus synthesized speech

Speaker-level volume analysis (Figure 8) confirms the systematic amplitude differences between natural and synthetic speech across all speakers in our corpus.

## C Evaluation Metrics

Our evaluation employs standard, well-established metrics from the speech processing and natural



Figure 8: Speaker-wise mean volume comparison between natural and synthesized speech

language processing domains:

$$\text{Perplexity} = \exp\big(\text{CrossEntropy}(p, q)\big), \quad (1)$$

$$F_1 = 2 \cdot \frac{\text{Precision} \times \text{Recall}}{\text{Precision} + \text{Recall}}, \quad (2)$$

$$\text{WER} = \frac{S + D + I}{N}, \quad (3)$$

$$\text{MAE} = \frac{1}{n} \sum_{i=1}^{n} |P_i - A_i|, \quad (4)$$

$$\text{RMSE} = \sqrt{\frac{1}{n} \sum_{i=1}^{n} (P_i - A_i)^2}, \quad (5)$$

$$\text{ARR} = \frac{\big|\{\text{ words aligned within } \tau \}\big|}{N}. \quad (6)$$

Here, $p$ and $q$ denote the true and predicted distributions (perplexity). In WER, $S$, $D$, and $I$ are substitutions, deletions, and insertions, and $N$ is the number of reference words. In MAE and RMSE, $n$ is the number of predictions, with $P_i$ and $A_i$ the predicted and actual values for instance $i$. For ARR (Alignment Recall Rate), $\tau$ is the temporal tolerance for correct alignment (e.g., $\pm 50$ ms). Unless otherwise specified, we report a macro-averaged ARR: the ratio is computed in each 15-second window and then averaged over all windows.

## D SSML Annotation Example

Figure 9 illustrates a representative example of our automated SSML annotation, demonstrating the integration of prosodic tags with natural text to

enable fine-grained control over synthetic speech parameters.

```xml
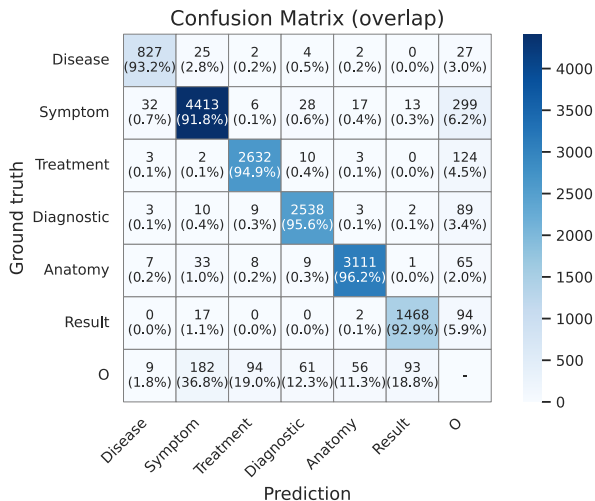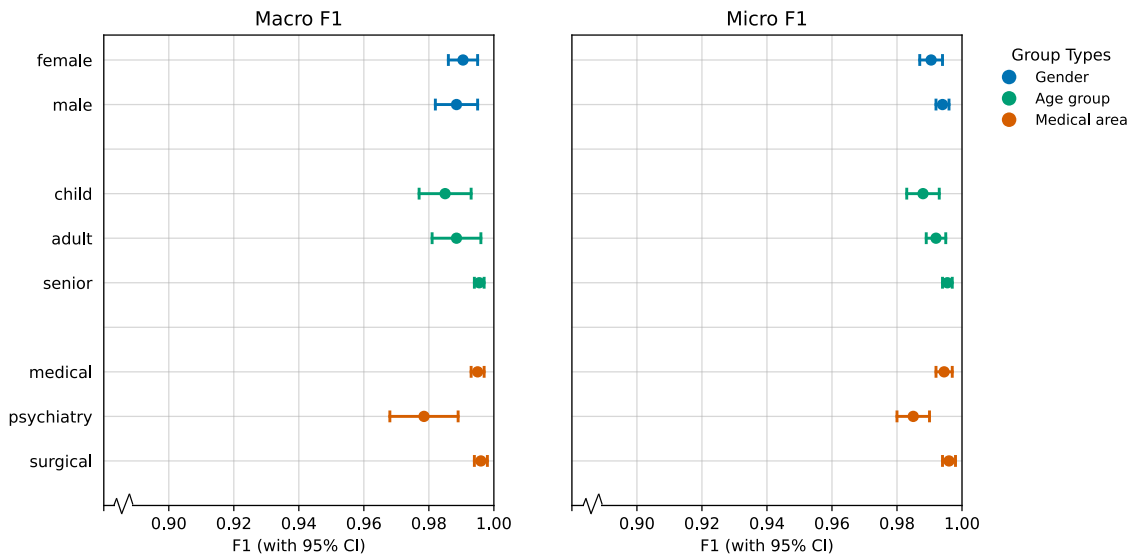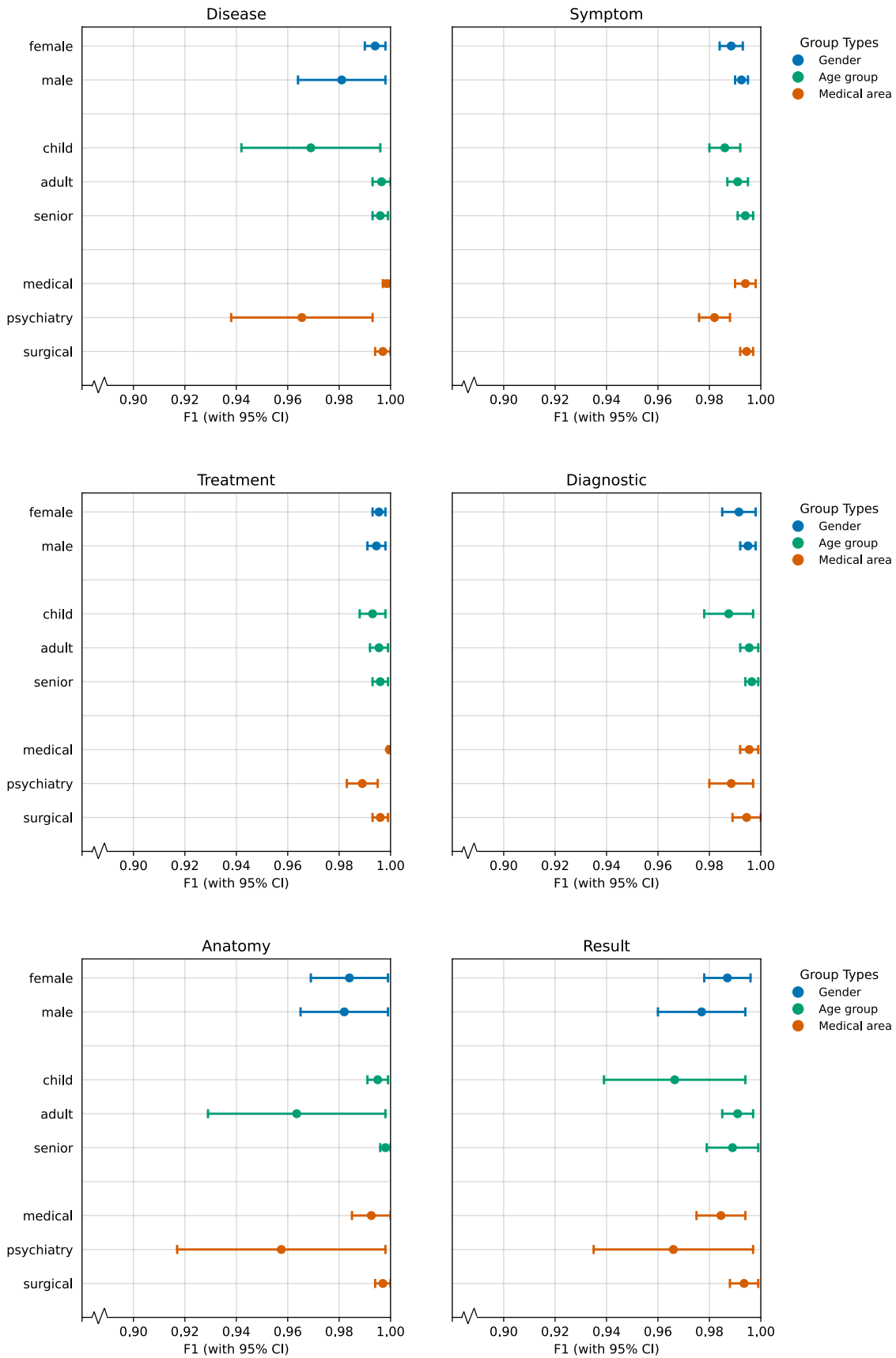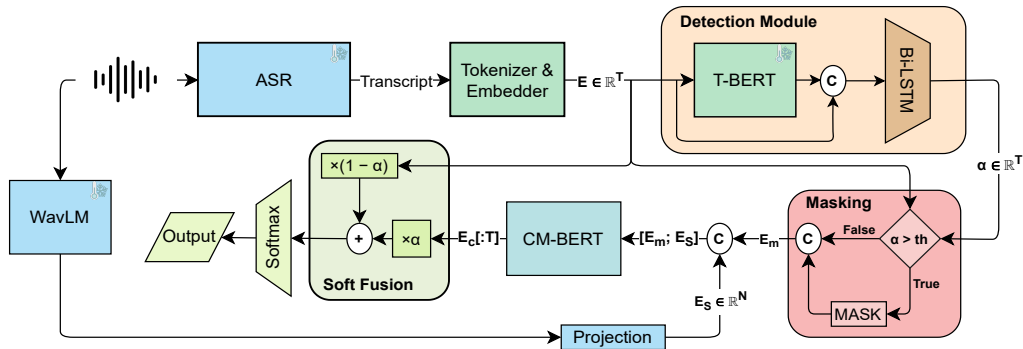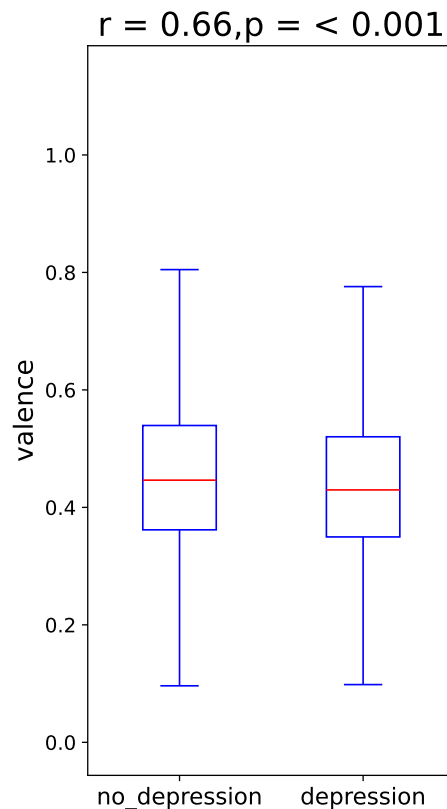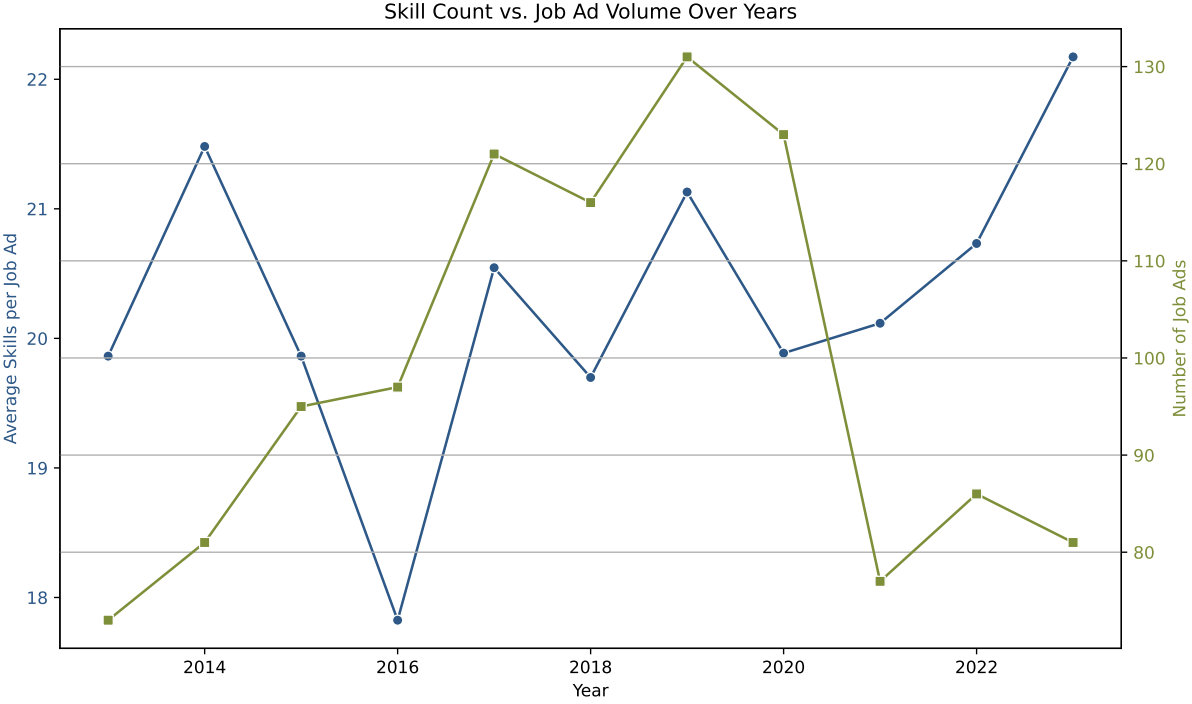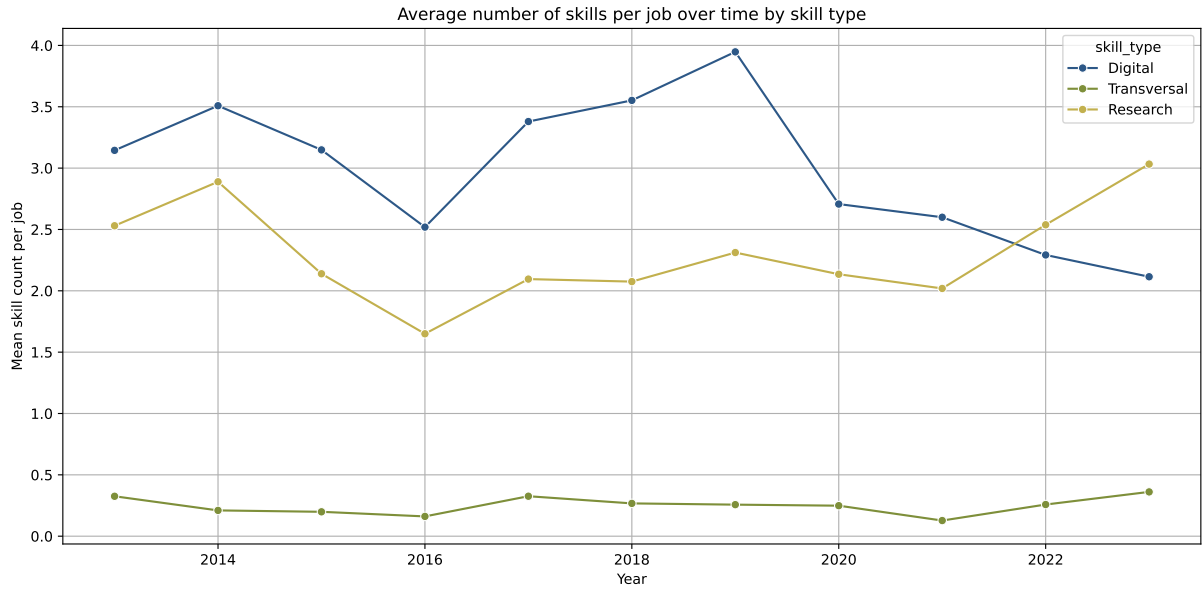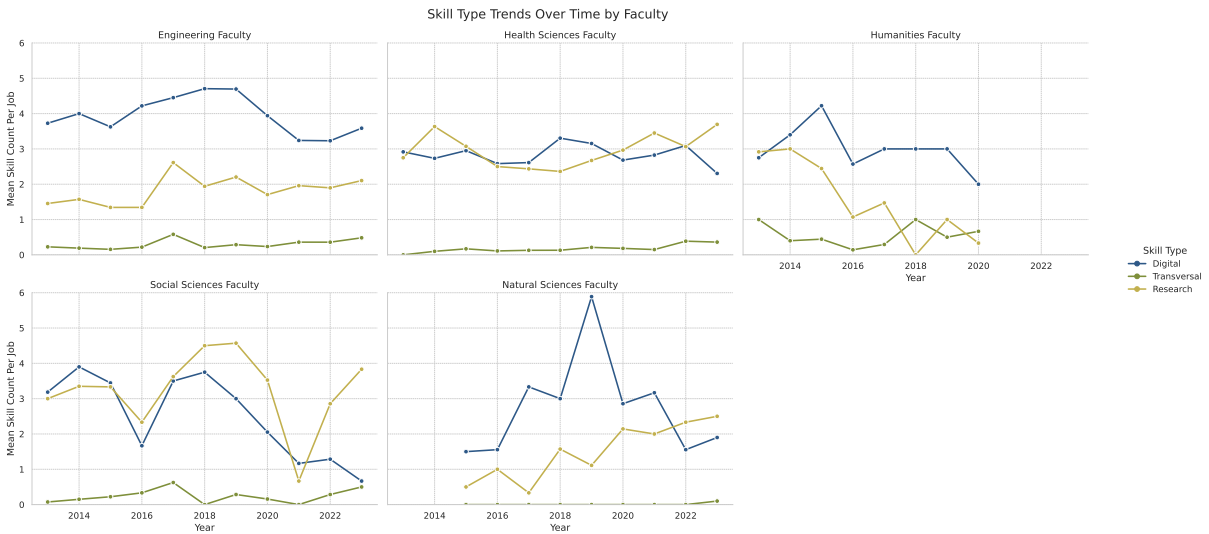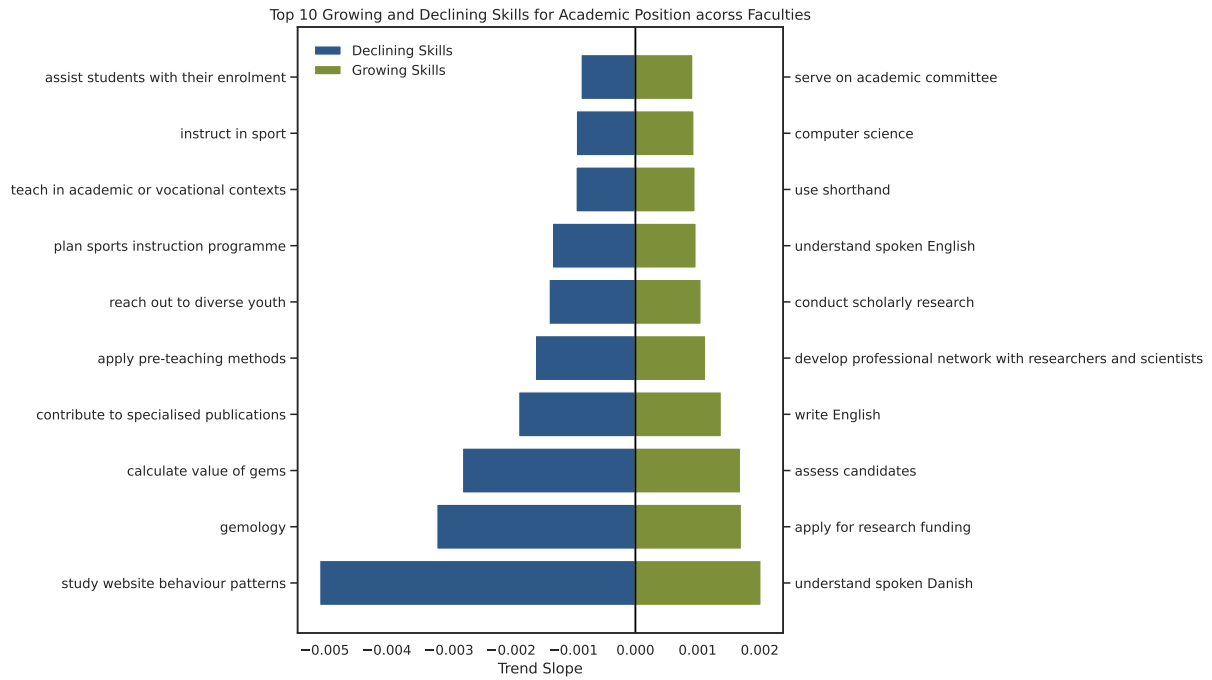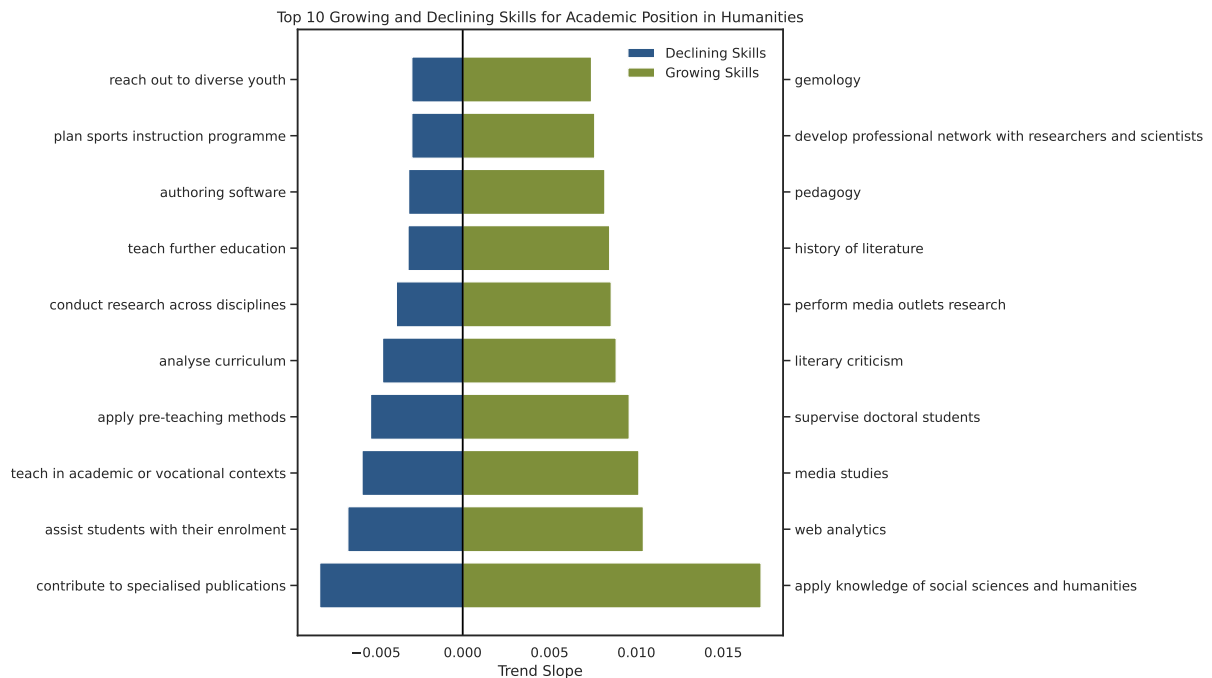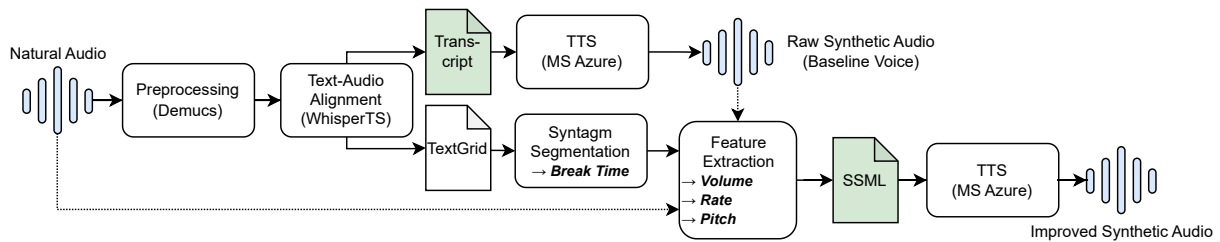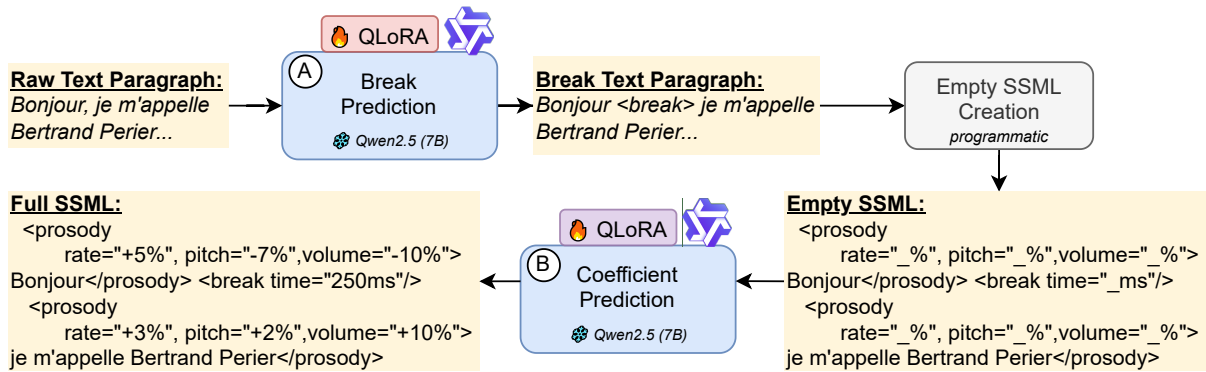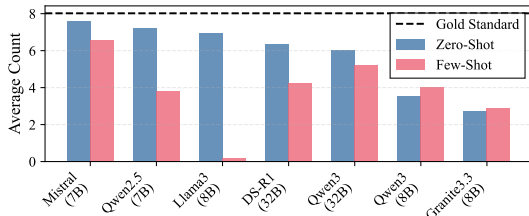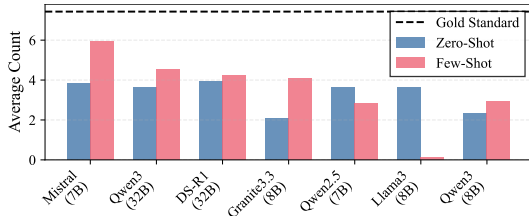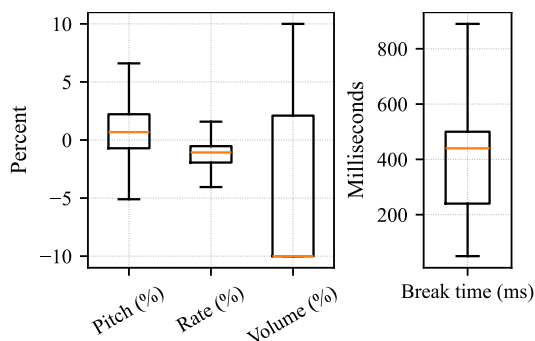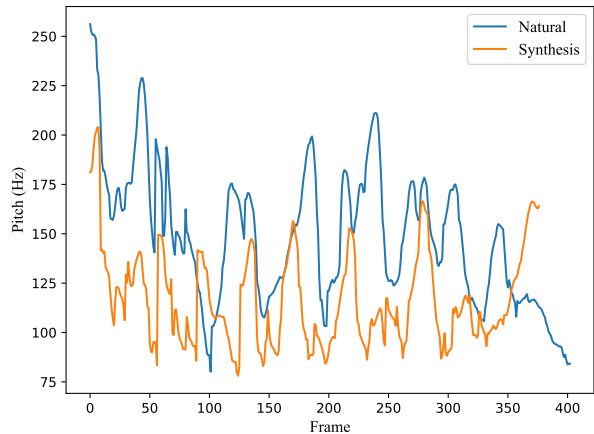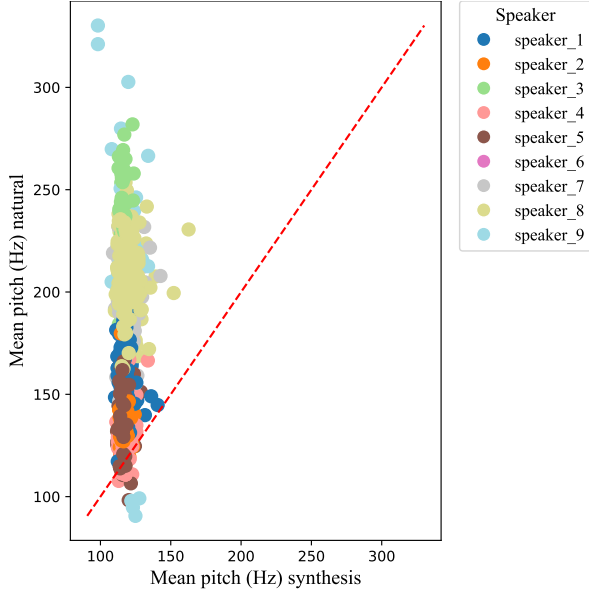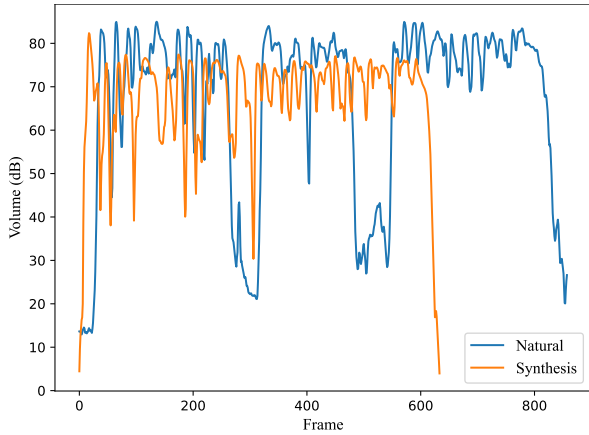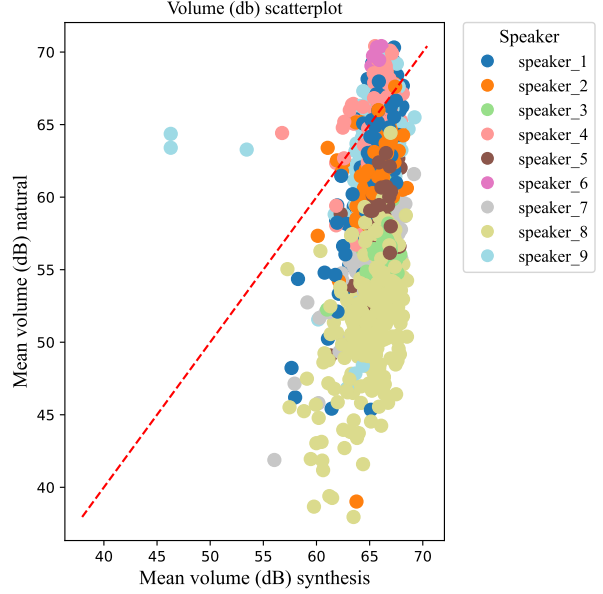<speak>
<prosody pitch="+2.01%" volume="+10.00%" rate="-3.10%">
Il y a dans la parole ce qu'on appelle la voix d'implication.</prosody>
<break time="500ms"/>
<prosody pitch="+2.73%" volume="+10.00%" rate="-2.18%">
Lorsque je vous parle actuellement,</prosody>
<break time="360ms"/>
<prosody pitch="+1.97%" volume="+10.00%" rate="-2.26%">
je fais un effort particulier pour moduler ma voix.</prosody>
</speak>
```

Figure 9: Example of text annotated with SSML prosodic tags generated by our pipeline

# Re-Representation in Sentential Relation Extraction with Sequence Routing Algorithm

**Ramazan Ali Bahrami** and **Ramin Yahyapour**
Georg-August-Universität Göttingen and GWDG , Göttingen, Germany
{ramazan.bahrami, ramin.yahyapour}@gwdg.de

## Abstract

Sentential relation extraction (RE) is an important task in natural language processing (NLP). In this paper we propose to do sentential RE with dynamic routing in capsules. We first show that the proposed approach outperform state of the art on common sentential relation extraction datasets Tacred, Tacredrev, Retacred, and Conll04. We then investigate potential reasons for its good performance on the mentioned datasets, and yet low performance on another similar, yet larger sentential RE dataset, Wikidata. As such, we identify noise in Wikidata labels as one of the reasons that can hinder performance. Additionally, we show associativity of better performance with better re-representation, a term from neuroscience referred to change of representation in human brain to improve the match at comparison time. As example, in the given analogous terms King:Queen::Man:Woman, at comparison time, and as a result of re-representation, the similarity between related head terms (King,Man), and tail terms (Queen,Woman) increases. As such, our observation show that our proposed model can do re-representation better than the vanilla model compared with. To that end, beside noise in the labels of the distantly supervised RE datasets, we propose re-representation as a challenge in sentential RE [1].

## 1 Introduction

Sentential relation extraction is about inferring the relation between two entities in a given sentence. Various sentential relation extraction datasets are constructed with distant supervision. Accordingly, it is assumed that if two entities are related in a knowledge base such as Wikidata, they are also related according to the sentence that contain them. As such, datasets often not only provide sentences with corresponding entities, but also additional details such as description, aliases, and types of entities. Accordingly, to improve performance, various



Figure 1: The need for re-representation in sentential relation extraction (Text from Wikipedia). At comparison time, the similarity between related terms increases. Note the entity types as manual label for re-representation. Note also the proportional word analogy constructed from the given example here, "Alan Turing": 1912 :: "C.F.GauSS":1777

works introduce complicated models that account for extra additional details. To that end, the results of incorporating additional details are not coherent. As example, while in some studies entity type has shown to improve performance (Bastos et al., 2021), in some others however, it has shown to degrade performance(Vashishth et al., 2018). As such, weather there are still room for improvement or if low performance are due to noise and error in the labels, is an open research question. In this paper, we aim to offer a deeper understanding of the task that can help in defining a better goal and objective for methods that incorporate the additional details into the sentential context using complex models. To do so, we propose an intuitive model that out perform state of the art on most dataset, and then investigate the reason for its better performance. To propose our approach, we build on works from neuroscience. To that end, inference such as sentential relation extraction in which a relation from one context, is mapped to a relation in another context is referred to as analogical reasoning (Gentner, 1983). In analogies, such as proportional word

---

[1] https://github.com/bahramiramazan/re-representation

analogies of the form King:Queen::Man:Woman, 'King' is related to 'Queen' as 'Man' is related to 'Woman', even though pairs (King,Man) or heads and (Queen,Woman) or tails are different. Similarly in sentential relation extraction, given that according to sentences $S_i$, for $i \in \{1, 2, ..., N\}$; the head entity, $e_i^h$ and tail entity, $e_i^t$ are similarly related, we have analogous terms of the form $e_i^h : e_i^t :: e_j^h : e_j^t$. To that end, studies in neuroscience suggest that, comparison as the foundation of any analogical reasoning, changes the representation of objects being compared (Gentner and Kurtz, 2006; Lorenza Saitta, 2013). As example in (Silliman and Kurtz, 2019) empirical evidence of representation, according to which people changes the representation of entities in order to improve the match at comparison time is documented. In other words, given that analogies are about partial similarity in different contexts (Hummel and Doumas, 2023), to map a relation from one context to a relation in another context, it is therefore needed to discard some information in both contexts (Lorenza Saitta, 2013) Figure 1. Accordingly, we found dynamic routing in Capsules network suitable for the task. Capsules were introduced first by (Hinton et al., 2011), and dynamic routing in capsules by (Sabour et al., 2017). They can be thought of neurons that output different features of processed entity. We test the proposed algorithm on relation extraction datasets Wikidata (Sorokin and Gurevych, 2017), Tacred (Zhang et al., 2017), Tacredrev (Alt et al., 2020), Retacred (Stoica et al., 2021), and Conll04 (Roth and Yih, 2004; Markus and Adrian, 2020). Our observations are summarized as follows:

- Our proposed approach improve state of the art scores on sentential relations extraction datasets Tacred, Tacredrev, Retacred, and Conll04.

- We estimate a significant error rate in labels of Wikidata, the dataset on which various studies try to improve model performance by incorporating the additional details through complex models.

- We show empirical evidence of re-representation and its associativity with better sentential RE performance in neural network.

## 2 Related Works

The use of extra additional details about entities such as entity type, aliases and description, and the way they are incorporated into the sentential context is one of the main theme of related works. As such, beside studies that addresses noise in the RE datasets, the other works deal with extra additional details and how to best incorporate them in the context. To that end, (Riedel et al., 2010) show that the vanilla distant supervised method used for generating sentential RE datasets, result in noisy labels, and proposes an improved version of the vanilla method, reducing error rate by 30%. Additionally, in studies related to variants of the common sentential RE dataset Tacred (Zhang et al., 2017), Tacredrev (Alt et al., 2020), and Retacred (Stoica et al., 2021); it is shown that after relabeling the noisy examples in Tacred, models improve performance by 8.0% (Tacredrev) and 14.3% (Retacred) of F1 score. Moreover, state of art performance for Tacred, and its variants, is proposed by (Zhou and Chen, 2022; Park and Kim, 2021). They show that incorporating abstract label of entities( entity types) improve model performance. Furthermore, (Sorokin and Gurevych, 2017) introduces Wikidata, a much larger dataset for sentential relation extraction based on the knowledge base wikidata (Table 1). To that end, to improve performance by enriching sentential context, in addition to entity type, (Nadgeri et al., 2021; Bastos et al., 2021) consider integrating other side information such as entity description, and aliases, through complex models such as graph neural network. Moreover, in (Vashishth et al., 2018), the use of entity types and relation alias information for improving performance is discussed.

## 3 Problem Formulation

We formulate the sentential RE based on the fundamental assumption that it is a type of analogical reasoning. To the end, in relation extraction datasets, according to sentences $S_i$, and $S_j$ we can construct proportional word analogies of the form $e_i^h : e_i^t :: e_j^h : e_j^t$, as we have King:Queen::Man:Woman. As such, a claim based on the studies from neuroscience is that our ideal proposed model shall do re-representation(Silliman and Kurtz, 2019). More commonly, given $X_i^h : X_i^t :: Y_j^h : Y_j^t$, re-representation can be viewed equivalent to a transformation $F$ such that according to some similarity measures $\psi$, when $X_i^h$ is related to $X_i^t$, as $Y_j^h$ is

Figure 2: Credit assignment in dynamic routing (Heinsen, 2022). The output here has 2 dimensions only, one for positivity, and the other for negativity. As example, given some sequence of vectors of depth h (25 here), sequence number n (number of tokens in the respective sentence), and dimension d (1024 here), and some configuration for the expected output (here depth=1, d=2, n=1), the dynamic routing algorithm works as credit assignment system. As such, projections of every feature in the input has limited credits at their disposal, and assigns it to the features in the output. Summing credits over all hidden states for positive feature, result in a value that is greater when the example is positive, and smaller otherwise(Example sentences are taken from Retacred).

related to $Y_j^t$ (positive examples) , we have :

$$\psi(F(X_i^h), F(Y_j^h)) \simeq 1$$
$$\psi(F(X_i^t), F(Y_j^t)) \simeq 1$$

and when $X_i^h$ is not related to $X_i^t$, as $Y_j^h$ is related to $Y_j^t$ (negative examples), we have:

$$\psi(F(X_i^h), F(Y_j^h)) \simeq -1$$
$$\psi(F(X_i^t), F(Y_j^t)) \simeq -1$$

As such, the similarity function $\psi$ returns 1 when terms come from positive examples and -1 otherwise. One common example for $\psi$ is cosine similarity between two given vectors. It is to note that, as in sentential RE, the sentences containing entities express the relation between entities, it is therefore needed that any change of representation be conditioned on the contextual sentence. As such, formally the task can be presented as the minimization problem below:

$$\min_{t \in [s,o]} \{\psi(F(X_i^t \mid S_i), F(X_j^t \mid S_j)) + (-1)^p\}_{i,j=1}^N$$

Here N stands for the number of instances in the dataset, p=1 when both entities come from the same relation, or from positive examples, and p=0 when entities come form negative examples. Additionally, $X_i^t$ stands for the embeddings of an entity or word, and $S_i$ is the embeddings for the contextual sentence. Moreover, $s$ stands for subject or head entity, and $o$ for object or tail entity.

## 4 Proposed Method

Our proposed model assumes an embedding model $\Omega$, and a transformation $F$ for obtaining the re-representations from the embeddings. Before giving a detailed description of our proposed method, we characterize it as follows.

1 Given some $X_i$, as tokens representing sentences, our transformation F obtains a single vector $x_{(out)}^{1d}$ of some dimension d for the joint representation of sentence containing entities $e_i^h$, and $e_i^t$, which are related to some relation $R = r_i$.

2 Instead of working with an explicit similarity function such as cosine similarity, our model is trained to maximize the following conditional probability:

$$P_\theta\Big(R = r_i \mid F\big(\Omega(X_i)\big)\Big)$$

With $\theta$ being the model parameters.

3 We show that maximizing the above conditional probability as we propose in this section will encourage explicit similarity, as was explained in Section 3.

Given that entities and sentences are sequence of vectors, our transformation $F$ can be such, given a sequence of vectors, it outputs a single vector $x_{(out)}^{1d}$. To that end, our proposed method for $F$ is dynamic routing in capsules. Capsules were introduced first by (Hinton et al., 2011), and dynamic

Figure 3: Over all architecture of our proposed model. We use different heads to classify the relation. On top is the Decoder, the head that translate from the sentence to entities post-fixed with relation id. Below it, are heads that implement routing as described in the paper (Heinsen, 2022). In that, H1( gray module in the middle) will identify positivity and negativity of examples. Under it is H3, the head that find the representation of the relation or sentence with entities marked as is shown in text on top of the diagram. Above H1 is H2, the head that calculate the joint representation of concerned entities. To the left is the pre-trained large language model(LLM), the backbone from which we obtain embeddings.

routing in capsules by (Sabour et al., 2017). In capsules network, neurons are grouped into capsules, each capsule representing some particular aspect or feature of the processed entity. Additionally, with dynamic routing, flow of data from a capsule in a layer to a capsule in the next layer depends not only on the weight matrix, but also on coefficients that itself depends on the data, also referred to as routing coefficients. The dynamic routing in capsules network is also called voting by agreement, as a capsule's vote is greater for capsules with which it agrees (Heinsen, 2022).

The dynamic routing algorithm used in this work (Heinsen, 2022), instead of voting by agreement, describe itself as credit assignment system Figure 2.

With that being said, to obtain the embeddings of the given sequence of words or sentence, we use some pre-trained model $\Omega$ such as bert_base_uncased (Devlin et al., 2019) or roberta-large (Zhuang et al., 2021):

$$X_{n(inp)}^{hd} = \Omega(X)$$

With $X$ being the tokenized sequence of words for

a sentence, and $n_{(inp)}$ being the number of tokens in our input sequence, $d$ being the embedding dimension, and h representing the number of hidden states in the pre-trained model.

To generate an output $x_{(out)}^{1d}$ as representation for the given sentence, the generated multi-dimensional matrix $X_{n(inp)}^{hd}$ is feed into the sequence routing algorithm (Heinsen, 2022). We call the routing module as the routing head. As such, the routing head used, is configured to convert a sequence of vectors of depth h, number of sequence n, representing a sentence, term or entity into a single vector of some dimension $d$. Additionally, for experiment, we also create some of the routing heads to do some specific predefined tasks, and evaluate if adding them to the main routing head can be of help. To that end, our main routing head is used for obtaining the representation of the sentence with some marking for the concerned entities as is shown in top of the Figure 3. Moreover, as in most datasets, a significant portions of the data are negative, we create an special routing head with two features, one representing positivity ( the relation between concerned entities is among our

| Dataset | Tacred | Tacredrev | Retacred | Conll04 | wikidata |
|---|---|---|---|---|---|
| No of Relations | 41 | 41 | 40 | 5 | 353 |
| No of Abstract Entities | 23 | 23 | 23 | 4 | 13533 |
| Train Size | 68,124 | 68,124 | 58465 | 1283 | 372,059 |
| Eval Size | 22,631 | 22,631 | 19,584 | - | - |
| Test Size | 15,509 | 15,509 | 13,418 | 422 | 360,334 |
| Negative Size | 79.5% | 79.5% | 79.5% | - | 29% |

Table 1: Statistics of datasets used in this work.

relation set) and another one representing negativity of examples. The working of the mentioned routing head is depicted in the Figure 2. A detailed description of routing heads and the baseline Decoder is depicted in the Figure 3. In practice, we experiment and compare performance of a single routing head, and all routing heads combined with the Decoder, referred here as collection of experts. Each head is characterized as follows:

1. H1: Obtains the representation for positivity or negativity. The output for this head has 2 dimension only, one representing positivity and another negativity Figure 2.

2. H2: This head learns the joint representation of head and tail entities.

3. H3: Is used to obtain the representation for the sentence containing concerned entities. H3 is the main routing head.

4. Decoder: We use a transformer based decoder for the baseline model, as is shown in the Figure 3. As in the example in the mentioned Figure, the decoder uses the last hidden state for the sentence as memory, and entities postfixed by the corresponding relation id as the target.

### 4.1 Optimization

Given the organization of our data into head and tail entities, $e_i^h$ and $e_i^t$, and the corresponding sentences $S_i$ for $i \in N$, and relation $r \in R$, with R being the relation set, and the embedding model $\Omega$, and the transformation $F$ based on dynamic routing, and an instance of data as follows:

$$\mathcal{D} = \{(e_i^h, e_i^t, S_i, r_i)\}_{i=1}^{N}$$

where $r_i \in R \quad and \quad i \in \{1, 2, 3, \ldots, N\}$

Where N is the dataset size. The transformation $F$ based on dynamic routing, learns a representation $x_i^d$, with some dimension d, such that the loss below is minimized:

$$\mathcal{L}(\theta) = -\sum_{i=1}^{N} \log P_\theta \Big( R = r_i \mid F\big(\Omega(X_i)\big) \Big)$$

Here $X_i$ being the tokens for the sentence $S_i$, and F is the routing head. In practice, by concatenating the outputs produced by different heads, we experiment if combining heads, also referred to as collection of experts Figure 3 may be of any help. Additionally, if Decoder is among selected heads, its loss will be added to the classifiers loss as in the Figure 3.

## 5 Experiments

### 5.1 Datasets

We test our model on several sentential relation extraction datasets. Specifically we test the proposed model on wikidata (Sorokin and Gurevych, 2017), Tacred (Zhang et al., 2017), Tacredrev (Alt et al., 2020), Retacred (Stoica et al., 2021), and Conll04 (Roth and Yih, 2004; Markus and Adrian, 2020). In all datasets, except Conll04, negative example make a significant portion of the examples. To that end, there are several factors to note about the datasets.

1. The ratio of positive and negative examples: Conll04 has no negative record, while wikidata has 22/29% of example as negative, and all Tacred variants have 79.5% of example as negative.

2. Number of entity types or manual abstract label of entities: All Tacred variants has 23 abstract labels for entities as and according to name entity recognition types in stanford

NER system (Zhang et al., 2017). Conll04 has only 4 different types of entities, and Wikidata has the highest number of abstract labels for entities 13,533.

3. Number of relations: From number of relation points of view, wikidata has 353 relation types, which is the largest among all datasets considered, while Conll04 has only 5 types of relations.

## 5.2 Different Configuration of the Sentence

The assumption for our proposed model is that, in order to do re-representation, dynamic routing can do feature selection . As such, we study our proposed model with different settings or configuration of the sentence. Accordingly, a given sentence can provide different level of details about the entities and their relations. As example, the sentence "<Mask> was getting married to <Mask>.", wherein the two concerned entities are masked, provide less details as when entities are not masked. Similarly, when entities are replaced by entity type, the level of details are less than the original sentence. This is important as in order to do re-representation, models need to do abstraction, and discard some unnecessary details. As such, manual abstract label of entities or entity types can perhaps make the job of RE models on some datasets easier. To that end, the following sentence configurations are used with markings as is shown in the Figure 3:

- Abstract: We replace surface form of the entity with entity type (abstract label of the entity). example: Germany or France is replaced by entity type Country.

- Mask: We replace surface form of the entity with the placeholder, 'MASK', in the sentence.

- Entities: We use only surface form of the entity as is.

- Mix: The entity type, or abstract label for entity and its surface form is used together with some marking. Example: "x was getting married to y." is transformed into : " [e11] + person * x [e12] was getting married to [e21] # person & y  [e22]."

| Config | Model | Retraced | Conll04 |
|--------|-------|----------|---------|
| Mix | H3 | **92.2(80.1)** | **100.0(100.0)** |
| | Decoder | 49.3(21.0) | 78.6(79.8) |
| Entities | H3 | **89.7(58.5)** | **84.1(84.7)** |
| | Decoder | 50.4(31.5) | 42.1(41.8) |
| MASK | H3 | **81.7(54.2)** | **80.1(79.3)** |
| | Decoder | – | – |
| Abstract | H3 | **75.2(48.5)** | **82.2(80.3)** |
| | Decoder | 29.1(13.0) | 61.8(63.7) |

Table 2: Comparative performance of the routing head H3, and transformer based Decoder on different configuration of sentence or information granularity. Recorded scores inside parenthesis are F1 Macro, and F1 Micro otherwise. The backbone model is roberta-large.

## 5.3 Experiment One: Comparative Performance on Different Information Granularity

We investigate performance achievable with our proposed model, and the transformer based Decoder on each sentence configuration described above. As each configuration of the sentence provide different level of details about entities and their relation, we refer to different sentence configuration as different information granularity. Accordingly, the relation between two entities in a sentence can be mostly predicted in all sentence configuration considered here; However, the best result by the proposed model is when the entity type is added to contextual sentence. For Decoder however, the best result changes across datasets considered Table 2. As such, on Retraced, Decoder's best score is when entity type is not added to the sentence( configuration "Entities"). However on Conll04, it is the other way around. Moreover, on Retraced, Decoder have relatively low scores, while on Conll04, our Decoder's score( 78.6 F1 Micro) is above state of the art ( with state of the art being 76.5, Tables 5, and Table 2).

## 5.4 Experiment Two: Entity Types as Manual Label for Re-Representation

Given that entity types increase the similarity as is expected for re-representation (a depicted example can be seen in the Figure 1), we can view entity types as manual label for re-representation. To that end, we extract entities from the respective sentences, and train the proposed model on the extracted entities and entity types. This help

us study entities and entity types in isolation. In doing so, we consider all sentence configurations (except Mask) as explained in the Section 5.2. As such, for Conll04, our proposed model exhibit same performance with configuration Abstract and configuration Mix Table 4. As such, we can conclude that, on some datasets, the manual label for re-representations or entity types (configuration Abstract), result in best performance. A possible explanation would be: when the entity types or manual label of re-representation can predict the relation, or given a relation the entity types for head and tail entities can be predicted, such as in Conll04, entity types alone (config Abstract) result in peak performance and less complexity Table 4.

## 5.5 Experiment Three: Performance on Varying Number of Entity Types

Does increasing the number of relation and entity types, or increased complexity for re-representation effect performance? To that end, we already observed the relatively good performance by Decoder on Conll04, the dataset with 4 entity types and 5 relations only Table 2. As such, we also evaluated the Decoder and the proposed model H3, on the smaller subset of Retacred, person-person, having only 6 relations and 1 entity type only. Additionally, after training on the full dataset, we recorded the performance on the same subset, person-person*. Accordingly, Decoder's performance is better when number of relation and entity types are smaller Table 3, as was noted for Conll04. As such, the experiment support the notion that transformer based Decoder changes performance across dataset presumably due to larger number of entity types, and relations. Unlike the Decoder, the proposed model exhibit relatively high performance across datasets, with different number of relation, and entity types.

| Dataset Subset | H3 | Decoder |
|---|---|---|
| Full | 92.2(80.1) | 49.3(21.0) |
| Person-Person | 93.0(82.2) | 72.6(60.4) |
| Person-Person* | 89.7(78.3) | 51.6(38.2) |

Table 3: Performance on varying number of relation and entity types (config mix).Values inside the parenthesis are F1 Macro, and F1 micro otherwise. Person-Person is the subset of Retacred having head and tail entity types as person only. It is the largest subset of Retacred categorized by head-tail entity types. Full is the entire dataset. Person-Person* is performance on the same subset, but by the model trained on the full Retacred.

| Metrics | Retacred | Conll04 |
|---|---|---|
| Mix | **71.7** | **100.0** |
| Entities | 71.3 | 48.1 |
| Abstract | 62.0 | **100.0** |

Table 4: RE using entities extracted from the sentence, and with routing head H2. The backbone model here is roberta-large. The reported values are F1 micro.

## 5.6 Experiment Four: Comparison with State of the Art

To compare with state of the art, we trained our proposed model on the mentioned datasets, and documented the result. The result is shown in the Table 5. Our observations show that our proposed model outperforms state of the art on 4 datasets. To that end, our routing head H3, with roberta-large as the backbone, keeps a relatively high performance on all datasets. It outperform state of the art on all dataset, except Wikidata. In the the Section 6.2 we show that noise is the main reason for the low performance on Wikidata. Moreover, despite the extra complexity that use of all heads, or expert heads, adds to our main model, we noticed little improvement. We therefore did not evaluate the expert head on Wikidata. Lastly, for our proposed model H3, the difference with bert-base-uncased (Devlin et al., 2019) and roberta-large (Zhuang et al., 2021) as the backbone is noticeable.

## 6 Observations

### 6.1 Re-Representation in Neural Network

As suggested initially, treating sentential RE as analogy, requires some form of re-representation to improve the match. To check if neural-network also does re-representation, using a subset from Retacred test set, we create positive and negative analogous entities of the form $e_i^h : e_i^t :: e_j^h : e_j^t$, for all $i, j \in \{1, 2, .., N\}$ such that the corresponding sentence $S_i$ and $S_j$ expresses the same relation between the corresponding entities in positive examples and different relation in negative examples. In doing so, we obtain the embedding for a given entity in the sentence, by feeding the sentence into the backbone model, and then slice the entity from the sentence embedding. We then calculate the cosine similarity, and pairwise euclidean distance between respective head, and tail entities in both positive and negative examples. We calculate the mentioned values across hidden states

(a) As can be seen, before training, the similarity between head and tail terms in positive (Heads +, Tails +) and negative (Heads -, Tails -) examples are barely distinguishable. However, after training, the model based on dynamic routing, does a good job of making head/tail terms more similar in positive examples, and dissimilar in negative examples.



(b) The distinction between positive and negative examples are barely distinguishable before training both for head terms (Heads +, and Heads -) and also for tail terms(Tails +, Tails -). However, after training, and that also specially for the model based on dynamic routing (*_route +, *_route -) the increase in the distance between head/tail terms in positive examples, are far less intense than in negative examples.

Figure 4: X-axis represent different hidden layers of the pre-trained LLM. Y-axis represent categories for which representation's similarity or distance was calculated. + represent positive analogous examples, and - represent negative analogous examples respectively. Heads and Tails are the related head and tail terms in proportional analogy. As example in king:queen::man:woman , (king, man) are head, whereas (queen, woman) are tail. We report the result of calculations obtained on representation after training with routing heads H3(Heads/Tails_route +/-) , and transformer based Decoder (Heads/Tails_decoder -/+). We also report the same before training (Heads/Tails +/-).

| Model | Tacred | Tacredrev | Retacred | Conll04 | Wikidata |
|---|---|---|---|---|---|
| Entity Marker (2022) | 74.6 | 83.2 | 91.1 | - | - |
| Curriculum Learning(2021) | 75.2 | - | 91.4 | - | - |
| REBEL (2021) | - | – | 90.4 | 76.5 | - |
| KGpool (2021) | - | - | - | - | **88.6** |
| RAG4RE (2024) | 86.6 | 88.3 | 73.3 | - | **-** |
| Ours bert ₍H3₎ | 84.8 (47.8) | 85.3 (49.7) | 89.4 (74.0) | 99.7(99.8) | 84.5 (32.0) |
| Ours bert ₍H1,H2,H3,Decoder₎ | **87.4** (48.3) | 88.7(50.9) | 88.7 (68.5) | **100.(100.)** | – |
| Ours Roberta ₍H3₎ | 87.1 **(61.1)** | **88.8 (64.2)** | **92.2(80.1)** | **100. (100.)** | 85.6 (32.9) |

Table 5: Our method's performance compared with state of the art. Best score is bold, state of the art is blue. Values inside the parenthesis are F1 Macro, and F1 micro otherwise. The configuration of sentence is Mix, and backbone is as indicated. We do not test all heads(H1,H2,H3,Decoder) for Wikidata as we found H3's performance to be already good on Wikidata's noisy labels.

of pre-trained backbone model, both after training with each training heads H3 and Decoder, and also before training, and then create a heat map as is shown in the Figure 4a. Accordingly, "Head +" , and "Head -" represent the cosine similarity between heads in positive and negative examples before training. As can be seen, the similarity is not much different between positive and negative examples. However, for the proposed model, after training, the similarity decreases significantly in negative examples, making the difference between positive and negative examples clearly noticeable ( specially in final layers of the backbone model). Similarly, the pairwise euclidean distance between positive and negative examples, shown in Figure 4b, after training are clearly distinguishable for the proposed model(Heads_route +, Head_route -) than it is for the vanilla Decoder (Heads_decoder +, Head_decoder -).

## 6.2 Noise in Wikidata's Labels

The tow Tacred variants (Tacredrev, Retacred) are very good attempts to improve data quality and reduce error rate in the Tacred. Each of these datasets improve model performance with 8.0% and 14.3% F1 score over the original Tacred respectively. In comparison to Tacred, wikidata has much larger and diverse types of relations. Its quality however has not gone a similar study. Instead, a significant attention has been given in improving model performance by incorporating extra additional details through complex models. As our model's performance is below state of the art on Wikidata, we were intrigued to have a look at examples in which our model disagree with labels from the dataset.

Not surprisingly though, we found out that a significant portion of errors are due to confusion in the dataset labels. As example, for instances which our model disagree with the dataset label, the labels seem random. More such examples, and statistics in Appendix B.2. We categorize all examples that our model disagree with labels from dataset in the appendix B.2, Table 6 , Table 7, and Table 8.

## 7 Limitations

Over all the dynamic routing proposed by (Heinsen, 2022) is efficient and scalable as is explained in the original paper. However using all routing heads as collection of experts increases the complexity n folds, where n is the number of routing heads in the model. However, the good news is that, perhaps a single H3 head can do a better job as is shown in the Table 5.

## 8 Conclusion and Future Research Directions

In this paper we improve sentential relation extraction performance on several benchmarks. Additionally we identify noise as one of the main cause for low performance on largest sentential RE benchmark Wikidata. Furthermore, we propose re-representation as one of the challenges of sentential RE models. Lastly, we show that sentential RE dataset may not be as much sentence dependent as expected B.1. For future research direction, we are planing to study word analogies of the form a:b::c:d, jointly with sentential RE datasets. Specifically, it would be interesting to see how much improvement can training sentential RE benchmarks bring to word analogy benchmarks.

## Ethics Statement

We did pay for the dataset Tacred. Other datasets are publicly available. In addition for the proposed algorithm being efficient, we tried to minimize our $CO_2$ footprint. As such, this work comply and adhere to ACL code of ethics.[2]

## References

Christoph Alt, Aleksandra Gabryszak, and Leonhard Hennig. 2020. TACRED revisited: A thorough evaluation of the TACRED relation extraction task. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 1558–1569, Online. Association for Computational Linguistics.

Anson Bastos, Abhishek Nadgeri, Kuldeep Singh, Isaiah Onando Mulang, Saeedeh Shekarpour, Johannes Hoffart, and Manohar Kaul. 2021. Recon: Relation extraction using knowledge graph context in a graph neural network. In *Proceedings of the Web Conference 2021*, WWW '21, page 1673–1685, New York, NY, USA. Association for Computing Machinery.

Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. 2019. Bert: Pre-training of deep bidirectional transformers for language understanding. Preprint, arXiv:1810.04805.

Sefika Efeoglu and Adrian Paschke. 2024. Retrieval-augmented generation-based relation extraction. Preprint, arXiv:2404.13397.

Dedre Gentner. 1983. Structure-mapping: A theoretical framework for analogy. *Cognitive Science*, 7(2):155–170.

Dedre Gentner and Kenneth Kurtz. 2006. Relations, objects, and the composition of analogies. *Cognitive science*, 30:609–42.

Franz A. Heinsen. 2022. An algorithm for routing vectors in sequences. Preprint, arXiv:2211.11754.

Geoffrey E. Hinton, Alex Krizhevsky, and Sida D. Wang. 2011. Transforming auto-encoders. In *Artificial Neural Networks and Machine Learning – ICANN 2011*, pages 44–51, Berlin, Heidelberg. Springer Berlin Heidelberg.

Pere-Lluís Huguet Cabot and Roberto Navigli. 2021. REBEL: Relation extraction by end-to-end language generation. In *Findings of the Association for Computational Linguistics: EMNLP 2021*, pages 2370–2381, Punta Cana, Dominican Republic. Association for Computational Linguistics.

John E. Hummel and Leonidas A. A. Doumas. 2023. Analogy and similarity. In Ron Sun, editor, *The Cambridge Handbook of Computational Cognitive Sciences*, 2 edition, Cambridge Handbooks in Psychology, pages 451–473. Cambridge University Press.

Jean-Daniel Zucker Lorenza Saitta. 2013. *Abstraction in Artificial Intelligence and Complex Systems*. Springer, New York, NY.

Eberts Markus and Ulges Adrian. 2020. *Span-Based Joint Entity and Relation Extraction with Transformer Pre-Training*. IOS Press.

Abhishek Nadgeri, Anson Bastos, Kuldeep Singh, Isaiah Onando Mulang', Johannes Hoffart, Saeedeh Shekarpour, and Vijay Saraswat. 2021. KGPool: Dynamic knowledge graph context selection for relation extraction. In *Findings of the Association for Computational Linguistics: ACL-IJCNLP 2021*, pages 535–548, Online. Association for Computational Linguistics.

Seongsik Park and Harksoo Kim. 2021. Improving sentence-level relation extraction through curriculum learning. Preprint, arXiv:2107.09332.

Sebastian Riedel, Limin Yao, and Andrew McCallum. 2010. Modeling relations and their mentions without labeled text. In *ECML/PKDD*.

Dan Roth and Wen-tau Yih. 2004. A linear programming formulation for global inference in natural language tasks. In *Proceedings of the Eighth Conference on Computational Natural Language Learning (CoNLL-2004) at HLT-NAACL 2004*, pages 1–8, Boston, Massachusetts, USA. Association for Computational Linguistics.

Sara Sabour, Nicholas Frosst, and Geoffrey E. Hinton. 2017. Dynamic routing between capsules. In *Proceedings of the 31st International Conference on Neural Information Processing Systems*, NIPS'17, page 3859–3869, Red Hook, NY, USA. Curran Associates Inc.

Daniel C. Silliman and Kenneth J. Kurtz. 2019. Evidence of analogical re-representation from a change detection task. *Cognition*, 190:128–136.

Daniil Sorokin and Iryna Gurevych. 2017. Context-aware representations for knowledge base relation extraction. In *Proceedings of the 2017 Conference on Empirical Methods in Natural Language Processing*, pages 1784–1789, Copenhagen, Denmark. Association for Computational Linguistics.

George Stoica, Emmanouil Antonios Platanios, and Barnabás Póczos. 2021. Re-tacred: Addressing shortcomings of the TACRED dataset. In *Thirty-Fifth AAAI Conference on Artificial Intelligence, AAAI 2021, Thirty-Third Conference on Innovative Applications of Artificial Intelligence, IAAI 2021, The Eleventh Symposium on Educational Advances in Artificial Intelligence, EAAI 2021, Virtual Event, February 2-9, 2021*, pages 13843–13850. AAAI Press.

---

[2]https://www.aclweb.org/portal/content/acl-code-ethics

Shikhar Vashishth, Rishabh Joshi, Sai Suman Prayaga, Chiranjib Bhattacharyya, and Partha Talukdar. 2018. RESIDE: Improving distantly-supervised neural relation extraction using side information. In *Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing*, pages 1257–1266, Brussels, Belgium. Association for Computational Linguistics.

Yuhao Zhang, Victor Zhong, Danqi Chen, Gabor Angeli, and Christopher D. Manning. 2017. Position-aware attention and supervised data improve slot filling. In *Proceedings of the 2017 Conference on Empirical Methods in Natural Language Processing (EMNLP 2017)*, pages 35–45.

Wenxuan Zhou and Muhao Chen. 2022. An improved baseline for sentence-level relation extraction. In *Proceedings of the 2nd Conference of the Asia-Pacific Chapter of the Association for Computational Linguistics and the 12th International Joint Conference on Natural Language Processing (Volume 2: Short Papers)*, pages 161–168, Online only. Association for Computational Linguistics.

Liu Zhuang, Lin Wayne, Shi Ya, and Zhao Jun. 2021. A robustly optimized BERT pre-training approach with post-training. In *Proceedings of the 20th Chinese National Conference on Computational Linguistics*, pages 1218–1227, Huhhot, China. Chinese Information Processing Society of China.

## A Training and reproducibility

For all routing heads we use the code from (Heinsen, 2022). Additionally, we use tokenizers from https://huggingface.co for bert-base-uncased, and roberta-large respectively. Our backbone models are too from https://huggingface.co. Furthermore, for routing head H3, we train it on datasets Tacred, Retacred, and Tacredrev with batch size 64, learning rate $10^{-5}$ , and on the dataset Wikidata with batch size 128, and similar learning rate as for Tacred and its variants. For collection of experts we use an smaller batch size of 24. For optimizer with use Adam from torch.optim. Moreover, we find hidden state of routing heads to have great influence on performance. To that end, for H3, w used hidden_d=256, and out_dimension=512. Moreover, we trained the proposed model for Tacred and its variants for 6 epochs, while we trained only for 1 epoch on Wikidata.

Another point to note is: In case of wikidata, when entities did not have an entity type ( or instance of as in the dataset), we checked the Wikidata knowledge base to retrieve parent class as entity type[3] . Furthermore, when entities had several values as "instance of" or parent class, again we query Wikidata to check if they have a common parent class, and used the parent class as the entity type, if not, the most common class was uses for entity type. Lastly, unless explicitly mentioned, all experiments are done with Roberta-Large as backbone.

## B Observation

### B.1 Are Sentential RE Datasets Truly Sentential?

To answer if relation between the entities, can be inferred without reading the sentence, and only be looking into entities, we trained and evaluated the proposed model on entities with configuration as was discussed for the sentence. The result for different configurations are recorded in the Table 4. Accordingly, most relation can be inferred without reading the concerned sentences.

### B.2 Noise in Wikidata's Labels

On examples which our model disagree with the dataset Wikidata, we found a pattern. Specifically, given a pair (p0-p*), where p0 is label('no relation') provided by the dataset, and p*(some relation other

than "no relation") predicted label, there is usually another category of predictions as (p*-p0). In both groups of examples, the probability that p * is true is similar, regardless of the label provided by the dataset. The group pairs, such as p0-p17 and p17-p0; show confusion caused as a result of incorrect labels. Some examples in Table 6.

For ease of understanding, we list Wikidata relation codes used in the table with corresponding labels as: • P131(located in the administrative territorial entity) • P17(country) • P47(shares border with) • P118(league) • P571(inception) • P47(shares border with) • P361(part of ) • P463(member of)

| Label-Prediction | probability of P* being True | count |
|---|---|---|
| P0-P17 | 80.6 | 4090 |
| P0-P131 | 90.0 | 4037 |
| P0-P47 | 60.0 | 3518 |
| P0-P118 | 70.0 | 2155 |
| P0-P571 | 50.0 | 1718 |
| P0-All | 70.0 | 29021 |
| P47-P0 | 60.0 | 12184 |
| P131-P0 | 80.0 | 4775 |
| P17-P0 | 70.0 | 4312 |
| P361-P0 | 60.0 | 2152 |
| P463-P0 | 70.0 | 1546 |
| All-P0 | 60.0 | 40155 |
| label!=prediction | - | 106534 |

Table 6: Top categories(sorted) on which model's predictions does not match with the label from benchmark. * represent a relation other than 'no relation'. The probability here is calculated by sampling 10 random example from each category, and then manually checking if p* holds.

[3] https://query.wikidata.org

| Label-Prediction | Sample |
|---|---|
| P0-P17<br>P0:No relation<br>P17:country | Los Dominicos is a metro station on Line 1 of the Santiago Metro in Santiago , Chile, and is also the eastern terminal of this line . |
| P0-P131<br>P0:No relation<br>P131: Located in the administrative territorial entity | Boechout is a railway station in Boechout , Antwerp , Belgium . |
| P0-P47<br>P0:No relation<br>P47:Shares border with | There are now approximately twenty restaurants in operation in Georgia , and about nine more in North Carolina , South Carolina , Florida , and Tennessee. |

Table 7: Random sample from $P0 - P^*$, where $p^*$ is any relation from relation set other than no relation, and p0 is no relation

| Label-Prediction | Sample |
|---|---|
| P6-P138<br>P6:head of government<br>P138:named after | She later served in the Blair ministry under Prime Minister Tony Blair in a number of roles , becoming Britains first female Foreign Secretary in 2006 . |
| P264-P136<br>P264:record label<br>P136:Genre | CD1 is the unofficial name of an untitled album by English industrial music band Throbbing Gristle , released in 1986 through record label Mute . |
| P1416-P102<br>P1416:affiliation<br>P102:member of political party | Other famous Solidarity activists such as [e11] Anna Walentynowicz Solidarity  activists such as Anna Walentynowicz , Zbigniew Romaszewski and Antoni Macierewicz have visited the Basilica as well . |

Table 8: Random sample from $p^* - p^*$, where $p^*$ is any relation from relation set other than no relation. Consider the first row in which both label and prediction is correct.

# Building an Ewe Language Dataset: Towards Enhancing Automatic Speech Recognition Technologies for Low-Resource Languages

Isaac Wiafe[1]    Akon Obu Ekpezu[2*]    Raynard Dodzi Helegah[1]
Fiifi Baffoe Payin Winful[1]    Elikem Doe Atsakpo[1]    Charles Nutrokpor[1]
Kafui Kwashie Solaga[1]
[1]Department of Computer Science, University of Ghana
[2]Department of Information Processing Science, University of Oulu, Finland
akon.ekpezu@oulu.fi

## Abstract

Automatic Speech Recognition (ASR) systems rely on large-scale, high-quality training datasets. However, low-resource languages, such as Ewe, remain underrepresented in the development of these systems. This study presents the development of a large-scale open-source speech dataset for Ewe, a Niger-Congo language spoken across Ghana, Togo, and Benin. Using supervised crowdsourcing, participants recorded descriptions of preselected culturally relevant images using a customized Android app. We collected 203,336 validated speech samples (1130 hours) from 1937 speakers, along with 107 hours of transcribed audio. To demonstrate the utility of the dataset for ASR, we fine-tuned Whisper base models, which were originally trained on Shona and Yoruba. The evaluation results suggest that both base models adapted well to Ewe and achieved a word error rate of 37%, an orthographic error rate of 45%, and a character error rate of 12%. A qualitative error analysis identified challenges including orthographic inconsistencies, morphological complexity, phonetic confusion, and dialectical variations. Thus, highlighting the need for dialect-sensitive and morphologically aware ASR modeling. The open-source release of this dataset provides a critical resource for advancing ASR research and linguistic preservation efforts for underrepresented African languages. Future work will explore self-supervised learning techniques to further improve performance using the unlabeled Ewe speech corpus.

## 1 Introduction

Africa's linguistic diversity poses major challenges for automatic speech recognition (ASR) because of the limited availability of high-quality open-source speech and text data for low-resource languages (LRLs). This is further exacerbated by the fact that over 1500 languages are endangered and may be lost by the end of the century (Bromham et al., 2022). Thus, prioritizing endangered languages via ASR development is crucial to preserving linguistic heritage, ensuring inclusivity (Chizzoni and Vietti, 2024), and preventing the loss of valuable cultural and historical knowledge (Jimerson et al., 2018).

Although natural language processing (NLP) has made significant progress particularly, ASR modeling in high-resource languages such as English, Mandarin, and Spanish, only a small fraction of the world's languages are supported by these technologies (Peterson et al., 2021). This is evident in the limited availability of annotated datasets, computational tools, and research funding for LRLs. Interestingly, this is not different in Ghana, which is a multilingual country with over 80 languages, yet all are LRLs. Existing multilingual speech datasets such as the Common Voice project (Ardila et al., 2020), African Speech Dataset (Olatunji et al., 2023), and GlobalPhone (Schultz, 2002) do not include the Ewe language.

Ewe is a Niger-Congo language spoken by approximately eight million people in Ghana and neighboring countries, Togo and Benin and neighboring countries. Yet it lacks the diverse corpora needed to develop ASR models and speech technologies. Although some studies (Antwi-Boasiako and Agyekum, 2022; Dei, 2024) have attempted to document and preserve some Ghanaian languages including Ewe, these initiatives are constrained by cost,

lack of expertise, and technological support. These languages are also limited by standardized spelling conventions, dialectal variations, unspecified orthographies, and potential code-switching. For instance, Ewe has different dialectal variants across different regions, primarily in terms of orthography and pronunciation. These variations stem from the fact that Ewe is spoken across multiple countries; Ghana (including different regions), Togo, and Benin, and each of these countries or regions have its own sociolinguistic influences. While its core grammatical structure remains consistent, dialectal differences manifest in phonetics, vocabulary, and spelling conventions (Sam and Agbloe, 2024). Additionally, digitization efforts are limited by the lack of a standardized Ewe digital keyboard that can be installed on computers. This complicates digitization and makes data collection and transcription more challenging. Other previous attempts to digitalize Ewe are limited to context such as religious texts (Resnik et al., 1999), which may introduce domain-specific biases during ASR model training.

## 2  Related Work

Current methods of speech data collection including sentence reading and uncontrolled crowdsourcing are expensive, time-consuming, and logistically complex. This makes large-scale dataset development in low-resource environments (LREs) challenging. Although existing speech data collection approaches have been demonstrated to be effective in some jurisdictions (Ragano et al., 2020; Panayotov et al., 2015) they may not be appropriate for collecting Ewe. For instance, the sentence reading approach utilized by studies (Ibrahim et al., 2022; Georgescu et al., 2020; Gutkin et al., 2020) may be ineffective for languages with limited standardized orthographies. Moreover, many indigenous speakers of Ewe may lack the functional literacy required to accurately read sentences written in Ewe. Given the linguistic complexity of Ewe, crafting sentence prompts that capture the full range of natural speech and spontaneous utterances would require considerable effort.

While (Callison-Burch and Dredze, 2010) utilized Amazon's Mechanical Turk for un-

controlled crowdsourcing speech data collection, this method may not be feasible in regions with limited digital literacy and internet access. Also, uncontrolled crowdsourcing (Ardila et al., 2020) often results in inconsistent recording conditions, varying audio quality, and a lack of standardized quality checks, which affects the reliability and usability of the resulting dataset. This necessitates the design of a more structured and contextually appropriate approach for collecting Ewe speech data.

Accordingly, this study seeks to use a scalable and cost-effective approach to collect, curate, and evaluate a large speech dataset for Ewe. Specifically, it aims to collect at least 1000 hours of spontaneous speech data and 100 hours of transcribed text in Ewe language. The dataset will be evaluated by training an ASR model in Ewe. This study is expected to make several key contributions to theory and practice. Firstly, it will address the critical data scarcity challenge by providing an open-source, large-scale, high-quality speech dataset for Ewe. This is expected to significantly expand available linguistic resources for ASR development. Also, by employing a scalable and cost-effective data collection approach, this study offers a replicable framework that can be adapted for other LRLs. This study will contribute to the development of ASR technology by leveraging existing ASR models such as Whisper Small to finetune and evaluate an ASR model for Ewe. Ultimately, this study seeks to advance linguistic preservation efforts, enhance digital inclusivity, and set a foundation for future advancements in speech technology for Niger-Congo languages.

## 3  Methods and Materials

### 3.1  Ewe Speech Data Collection Pipeline

Ewe is linguistically complex. Hence, although it has a simple grammatical structure that makes it easy to decompose polysyllabic words into monosyllabic roots, it is characterized by unique phonological, morphological, and syntactic features. It is a tonal language with three main tones (high, mid, and low), that are used to distinguish the meaning of words. Hence, a phonetic structure may have dif-

ferent meanings depending on the tone used (e.g., "to" means "mountain" in one tone and "ear" in another). This makes it a challenge for speech recognition when compared to non-tonal languages such as English. Ewe also has a complex morphophonemic process (vowel harmony and nasalization) which affects the pronunciation of words depending on their syntactic environment. It is characterized by significant dialectal variations, where there are differences in pronunciation, vocabulary, and grammar based on regions (Sam and Agbloe, 2024). These variations make the development of a standardized speech recognition system a challenge. Nonetheless, the development of an Ewe dataset would augment ASR research and provide opportunities to develop technologies for over eight million speakers to support education, healthcare, and government services.

Existing speech data collection approaches are not well-suited for a LRL such as Ewe. Thus, a more structured and contextually appropriate approach for collecting Ewe speech data is necessary. To determine the most efficient approach, a focus group discussion with both functionally and non-functionally literate participants was conducted. The discussion revealed that since this study seeks to ensure a diverse representation of the Ewe language, as well as capture all possible complexities and scenarios of the language, then sentence reading would be impractical. This is because most of the study participants would be unable to read Ewe. Thus, image descriptions were considered the most suitable approach for collecting speech data in Ewe. This approach will facilitate the collection of a diverse range of spoken words in the form of sentences and also address the challenges of performing sentence segmentation of audio data manually (Uliniansyah et al., 2016).

Over 8000 images were initially extracted from online sources including Pinterest and Google images. Out of which a subset of 1000 images cutting across 50 categories (such as Sanitation, Tourism, Weather, Technology, Automobile, Security, TransportatioRobbien, Architecture, Fashion, Food, Trading, Hospitality, Lifestyle, Health/Medicine, Agriculture, through Entertainment, Arts/crafts, Science, Mining, Education, Governance, Leisure,

Home/Housing, Religion, Engineering, Accidents, Sports, Culture, Family, and Nature) were selected during the focus group discussion. Selected images were required to be easily describable in at least three different ways between 15 seconds and 30 seconds. In addition, images were screened to ensure they were devoid of any nudity or profanity. Context specificity was another consideration. This was to ensure that the selected images were culturally and linguistically relevant for the native speakers of Ewe. The images were uploaded onto an Android mobile app (UGSpeechData) that was developed to collect the data. The images alongside the URLs were integrated into the app's image database. The app was designed to operate on-device and with/without the Internet. Figure 1 shows the data pipeline from image selection to data finalization.

## 3.2 Study Participants and Speech Data Collection

Almost 2000 volunteers from diverse Ewe-speaking regions were recruited using convenience sampling and snowballing. Participants signed up and were trained to use the app to record image descriptions following a set of predefined rules. They signed the consent form and provided relevant demographics, including their age, gender, and recording environment. In addition, the app retrieved the device's name and the recording timestamp. Subsequently, all this data was stored in a file linked with the audio files.

Participants were required to describe the selected 1000 images in Ewe. Each image was limited to a single recording by a specific participant, and the app would only allow recording to start when there was little or no background noise. Participants could save, replay, and delete their recordings. However, the app was designed to only permit an audio file to be saved if it was between 15 and 30 seconds; if there was less than a three-second pause during the description; and if there were no excessive speech mannerisms/fillers in the description.

Furthermore, to ensure that the recruited participants could speak Ewe fluently, they were initially assigned 10 images and were required to record descriptions of the 10 images

Figure 1: Data Collection Pipeline Using Image Prompts



Figure 2: Distribution of audio clip duration in the dataset

to check for adherence to the rules of recording, language fluency, and audio quality. Restrictions to continue recording audio descriptions of the remaining 990 images were removed by the authors if at least 8 out of the 10 descriptions were validated and accepted. An audio file was valid and accepted if: the image description was in Ewe, there was no conflicting background sound in the recording, the audio was naturally audible, the description matched the displayed image and did not contain excessive use of English words, or filler words. Participants with less than 8 accepted recorded audio files were blocked and compensated but could no longer participate in the study. Out of the 2000 participants who were initially recruited, 1905 including 1076 males and 816 females passed the pre-selection phase. Their ages ranged between 18 and 74 years old with a majority between 18 and 45. See Table 1 for a summary of the participant's demographics and the number of audio files. A total of 203,391 audio samples, equivalent to 1,198 hours were recorded. Although participants were required to provide audio descriptions of 1000 images, they were at liberty to stop the recordings at any point. The audio durations range from 15 to 30 seconds, with most clips concentrated between 15s and 20s, and a gradual decline in the number of longer clips from 21s to 30s. Figure 2 shows the distribution of audio duration in seconds.

## 3.3 Audio Validation and Transcription

Following the collection of speech samples, thirty participants who recorded were reassigned to validate the audio based on the predefined rules. They were further trained on the stringent validation rules specified earlier (see Section 2.1). Out of the 203,391

| Gender | No. of recorders | No. of audio files recorded |
|---|---|---|
| Male | 1076 | 121,116 |
| Female | 816 | 81,684 |
| Other | 13 | 591 |
| **Total** | **1905** | **203,391** |

| Age range | No. of recorders | No. of audio files recorded |
|---|---|---|
| 18–25 | 751 | 56,361 |
| 26–35 | 606 | 84,544 |
| 36–45 | 287 | 32,676 |
| 46–55 | 149 | 17,175 |
| 56–65 | 71 | 9,796 |
| 66–75 | 30 | 2,613 |
| Unspecified | 11 | 226 |

Table 1: Distribution of participants' demographics across audio recordings.

(1198 hours) collected speech samples, 203,336 (1130 hours) speech samples passed the quality checks and formed the Ewe speech dataset. Furthermore, twenty linguists were trained to transcribe the validated speech samples using a structured workflow. We sought to transcribe at least 100 hours out of the 1130 hours of validated speech samples. A maximum of 240 audio files were randomly assigned to a transcriber every 48 hours. Each file was transcribed by two linguists and in situations where there are conflicts in the transcription, the audio will be passed on to a third linguist for conflict resolution. To facilitate transcription, a custom Ewe keyboard was developed to incorporate diacritics and special characters essential to the language. The keyboard utilizes the standard QWERTY keyboard layout and incorporates all special characters to support the Ewe orthography (i.e., including diacritics and tonal marks). The Ewe alphabet consists of 30 characters including the 26 letters of the English alphabet, excluding c, j, and q which were replaced by ɔ, ɣ, and ʃ respectively. In addition to the standard alphabet, the Ewe keyboard includes these special characters: ɖ, ŋ, ɛ, ɔ, ɣ, ʋ.

## 4 The Ewe Speech Dataset and Automatic Speech Recognition Experiment

The generated dataset consists of 203,336 (1,130 hours) validated audio speech samples, along with 19,152 (106.4 hours) of transcribed

audio containing 31,756 unique words. Each audio file is between 15 and 30 seconds long.

Audio speech samples were received from participants in two regions of Ghana: Greater Accra and the Volta Region. Within the Volta Region, data was collected from eight towns, namely Anloga, Keta, Peki, Akatsi, Ho, Juapong, Kpando, and Sogakope. The recordings were done in different environments, but the majority were done outdoors. Specifically, 118,193 recordings were done outdoors, 74,169 were done indoors, 2,465 were done in offices, 66 were done in studios, 144 in buses, and 6,755 in unspecified environments. The dataset is open-source and available at GitHub and Science Databank (Wiafe et al., 2025). See Table 2 for a summary of the Ewe dataset. Next, using the transcribed audio recordings, we test the suitability of the generated speech corpus for automatic speech recognition and also conduct a qualitative error analysis of the predicted transcriptions.

graphicx

### 4.1 Data Preparation, Fine-tuning and Evaluation

The initial dataset used for modeling comprised 106.4 hours of transcribed audio, encompassing 19,152 audio files. To ensure data quality and to eliminate potentially invalid entries, audio files that exceeded 30 seconds in duration and transcriptions containing fewer than 10 characters were excluded. Following this refinement, the final dataset consisted

| Gender | Total no of audio files | Equivalent in hours | Outdoors | Indoors | Other | Office | Car | Studio | Bus |
|---|---|---|---|---|---|---|---|---|---|
| **Total no. of audio files by environment** | | | | | | | | | |
| Male | 121 116 | 673.98 | 68 300 | 46 162 | 3 438 | 1 850 | 931 | 412 | 23 |
| Female | 81 684 | 453.80 | 49 851 | 27 513 | 3 317 | 615 | 17 | 250 | 121 |
| Other | 536 | 2.98 | 42 | 494 | 0 | 0 | 0 | 0 | 0 |
| **Totals** | **203 336** | **1 130.76** | **118 193** | **74 169** | **6 755** | **2 465** | **948** | **662** | **144** |
| **Summary of Transcribed Files** | | | | | | | | | |
| Gender | Total no of audio files | Equivalent in hours | Outdoor | Indoor | Other | Office | Car | Studio | Bus |
| Male | 10 870 | 60.39 | 4 372 | 5 729 | 398 | 210 | 140 | 21 | 0 |
| Female | 8 282 | 46.01 | 4 929 | 3 107 | 114 | 80 | 0 | 52 | 0 |
| **Totals** | **19 152** | **106.4** | **9 301** | **8 836** | **512** | **290** | **140** | **73** | **0** |

Table 2: Summary of the dataset (validated and transcribed)

of 19,149 audio files with a sampling rate of 16kHz. The dataset was partitioned into training sets (13,382 files, 70%), test sets (3,847 files, 20%), development sets (1,535 files, 8%), and validation sets (385 files, 2%). In terms of speaker distribution, the training set included 163 unique speakers, while the test, development, and validation sets contained 137, 130, and 109 unique speakers, respectively.

We selected the Whisper Yoruba and Shona base models (Radford et al., 2022) as base models due to the linguistic similarities between Yoruba and Ewe. Both languages share a similar writing system, are tonal with three tone levels, and exhibit some lexical overlap. For example, "mouth" is enu in Yoruba and enu/nu in Ewe, and "father" is baba in Yoruba and papa in Ewe. We fine-tuned both base models on the prepared dataset using Google Colab with NVIDIA A100 GPUs.

## 4.2 Training setup

The model was fine-tuned with the following hyperparameters: a batch size per-device of 16, gradient accumulation steps of 1, and a learning rate of 1e-5. We used the AdamW optimizer and applied a constant_with_warmup learning rate scheduler with 50 warm-up steps. Mixed precision training (fp16) and gradient checkpoint were enabled to reduce memory usage. The training process consisted of 2400 steps, and we evaluated the model's performance every 400 steps using the Word Error Rate (WER), arthographic error rate (OER), and character error rate (CER) as the pri-

mary metrics. These are widely used metrics for evaluating ASR performance (Fatehi et al., 2025; Mensah et al., 2025). Table 3 summarizes the results of the training loss, validation loss, OER, CER, and WER achieved at each evaluation checkpoint for the Shona and Yoruba base model. It was observed that both models exhibited similar performance trends across metrics. While training loss consistently decreased throughout, validation loss began to plateau after approximately 1600 steps. The lowest error rates across all metrics were recorded at 2000 training steps, with the Yoruba base model achieving an OER of 44.98%, WER of 37.12%, and CER of 12.43%, and the Shona base model achieving an OER of 45.11%, WER of 37.17%, and CER of 12.50%. Beyond this point, error rates showed slight increases, suggesting possible overfitting. Consequently, the 2000-step checkpoint was selected as the best-performing model for Ewe ASR. These results suggest that both base models adapt well to Ewe, with the Yoruba base model slightly outperforming the Shona model on all error metrics. Table 4 shows sample transcriptions predicted by the final model and the corresponding original text using the validation set. Irrespective of the relatively high error rates, the model was observed to make intelligible transcriptions.

## 4.3 Qualitative Error Analysis on the Predicted Transcriptions

Although it may be argued that the WER of 37% is high, (Chizzoni and Vietti, 2024) posit

| Step | Training Loss | Validation Loss | OER (Shona/Yoruba) | WER (Shona/Yoruba) | CER (Shona/Yoruba) |
|---|---|---|---|---|---|
| 400 | 0.50 | 0.58 | 52.37/51.88 | 44.57/44.15 | 15.05/14.89 |
| 800 | 0.48 | 0.52 | 48.49/48.65 | 40.52/40.66 | 13.69/13.75 |
| 1200 | 0.38 | 0.49 | 47.10/46.80 | 38.72/38.46 | 13.22/13.03 |
| 1600 | 0.36 | 0.48 | 46.08/45.92 | 37.86/37.71 | 12.83/12.70 |
| 2000 | 0.31 | 0.48 | 45.11/44.98 | 37.17/37.12 | 12.50/12.43 |
| 2400 | 0.31 | 0.47 | 45.56/45.43 | 37.58/37.48 | 12.86/12.97 |

Table 3: Model performance for the Shona and Yoruba base model

that the CER is a better evaluation metric in instances where the base model was not trained on the LRL data in question. Regardless, a qualitative error analysis was conducted to understand factors contributing to the relatively high error rates (see Table 3). Although the model generally produced intelligible transcriptions, several recurring challenges were identified across orthographic, linguistic, and acoustic dimensions.

1) **Orthographic inconsistencies**

   a) **Non-standard spelling of English loanwords.** Because most loanwords lack fixed Ewe spellings, transcribers wrote them phonetically. *Example:* "machine" appeared as **masini** or **mashini**.

   b) **Dialectal vs. formal spellings.** Mixing of Southern-Ewe forms with the formal standard produced mismatches. *Example:* model: **yi nye** (Southern) vs. reference: **si nye** (formal).

2) **Morphological challenges** — the model sometimes mis-segmented Ewe's agglutinative morphemes.

   a) **Reference:** ...enye nugometsi kpakple agbalē gbadza aɖe...
   **Prediction:** ...enye nugo me tsi kpakple agbalē gbadza aɖe...
   **Error:** nugometsi → nugo me tsi

   b) **Reference:** Devia ɖewo tsi atsitre…
   **Prediction:** Devia ɖewo tsatsitre…
   **Error:** tsi atsitre → tsatsitre

   c) **Reference:** exɔtudzikpɔla aɖe le wo gbɔ

   **Prediction:** exɔ tu dzikpɔla aɖe le wo gbɔ
   **Error:** exɔtudzikpɔla → exɔ tu dzikpɔla

3) **Phonetic confusions** — substitutions between phonetically similar consonants, especially affricates vs. stops.

   a) /dz/ ↔ /z/
   **Reference:** Dzo bi teƒe sia
   **Prediction:** Zo bi teƒe sia
   **Error:** Dzo → Zo

   b) /dz/ ↔ /d/
   **Reference:** Nufiala dzidzim be ya fia nu
   **Prediction:** Nufiala didim be ya fia nu
   **Error:** dzidzim → didim

4) **Dialectal pronunciation variation** Ewe exhibits major dialectal differences. The model often defaulted to Southern-Ewe pronunciations, causing mismatches when the reference used another variety. *Example:*
   **Reference:** Yevuwo wonye (standard Ewe)
   **Prediction:** Yewuwo wonyo (Ewe-Dome dialect)
   **Error:** wonye → wonyo

5) **Mistranscription with insertion/substitution** Rare but notable cases where acoustically ambiguous segments led to entirely different words: *Example:*
   **Reference:** Buno aɖɛ le suku
   **Prediction:** Nubuno aɖɛ le suku
   **Error:** Buno → Nubuno

| Original Text | Predicted Text |
| --- | --- |
| Ŋutsu eᵗɔwo le mashinidɔwɔfe le dɔ wɔm, wo dometɔ eve tɔ ɖe gakpo gā lɔbɔ aɖe yi le dzi ŋu le ŋku lem ɖe eŋu. Ɖeka tɔ ɖe adzɔge le wo kpɔm. Wodo awu amadede si nye orange eye woɖɔ dɔwokukuwo hā. | Ŋutsu eᵗɔwo le machinidɔwɔfe le dɔ wɔm. Wo dometɔ eve tɔ ɖe gakpo gā lɔbɔ aɖe yi le dzi ŋu le ŋku lém ɖe eŋu. Ɖeka tɔ ɖe adzɔge le wo kpɔm. Wodo awu amadede yi nye ɔɖɔɛndzi eye woɖɔ dɔwɔ kukuwo hā. |
| *Three men are at work in a machine shop, two of them are standing on a big, long steel stick that is above and staring at it. One stood at a distance watching them. They wore orange outfits and helmets.* | *Three men are at work in a machine shop, two of them are standing on a big, long steel stick that is above and staring at it. One stood at a distance watching them. They wore orange outfits and helmets.* <br> WER=39%, CER=11%, Cosine Similarity=95% |
| Woɖo kpᴵɔ aᵗɔ ɖe xɔ me. Amewo nɔ kpᴵɔ ŋu hamehame. Ame bubu aɖewo nɔ wo ŋgɔ eye wonɔ nu ƒom na wo. Ame siwo le kplɔ ŋu la ɖo to hele amea ƒe nu ƒom sem. | Woɖo kplɔ aᵗɔ ɖe xɔ me. Amewo nɔ kplɔ ŋu hamehame. Ame bubu aɖe nɔ wo ŋgɔ eye wonɔ nu ƒom na wo. Ame siwo le kplɔ ŋu la ɖo to hele amea ƒe nu ƒom sem. |
| *They set up five tables in a room. There were all kinds of people around the table. There were others in front of them and talking to them. The people at the table were quiet and listening the talk.* | *They set up five tables in a room. There were all kinds of people around the table. There were others in front of them and talking to them. The people at the table were quiet and listening the talk.* <br> WER=18%, CER=4.5%, Cosine Similarity=98% |

Table 4: Sample of ground truth vs. predicted transcriptions

## 5 Discussion

The performance of the fine-tuned Whisper Yoruba model on the Ewe dataset, achieving a word WER of 37% and CER of 12% is consistent with expectations for low-resource environments (LREs). Previous studies (Fatehi et al., 2025), have shown that automatic speech recognition (ASR) is dependent on the volume and quality of training data. The Common Voice project (Ardila et al., 2020) demonstrated that while community-driven data collection efforts help address issues of limited labeled speech date, achieving low error rates remains challenging without significant resources for transcription standardization and quality control. In high-resource environments (HREs), models achieve significantly lower error rates of less than 10% because they are trained on tens of thousands of hours of annotated speech (Baevski et al., 2020). However, in LREs such as the Ewe language, even with 106 hours of transcribed data and dili-

gent data collection efforts, the model's performance may have been constrained by the relatively small labeled data, dialectal variation, and orthographic inconsistencies. Besacier et al. (2014) argue that irrespective of advanced modeling techniques, there is an elevated risk of error rates, particularly for tonal and morphologically rich languages that lack large, domain-specific corpora. Dialectal variation and phonetic diversity, as observed in this study for Ewe, introduce substantial complexity. Dialectal shifts across regions (Ghana, Togo, Benin) result in pronunciation and vocabulary differences that pose challenges for ASR systems trained on limited samples. Orthographic inconsistency further exacerbates model error rates, as shown by (Kim et al., 2025) who highlighted the difficulties of building reliable models for languages with non-standardized or emerging writing systems. More specifically, the qualitative error analysis showed that orthographic inconsistencies, particularly with English loanwords

and dialectal variations, introduced ambiguities that affected transcription accuracy. It was observed that the morphological complexity of the Ewe language, especially its agglutinative nature, led to frequent segmentation and merging errors. Additionally, phonetic confusion between similar sounds (e.g., /dz/ vs. /z/) and dialectal variations in pronunciation may have compounded the ASR model's challenges. Despite leveraging transfer learning from a linguistically related language (Yoruba), the results show that adaptation alone cannot fully resolve the dialectical diversity or phonetic complexity intrinsic to Ewe.

## 6 Conclusion

This study introduced a large-scale, validated speech corpus for the Ewe language, comprising 1,130 hours of audio recordings and 106 hours of transcriptions. By employing an innovative image-based prompting method and controlled crowdsourced data collection strategy, this study provides a linguistic resource for advancing ASR development in LRE. Fine-tuning experiments with the Whisper Yoruba model demonstrated the dataset's utility while also highlighting persistent challenges posed by dialectal variation, orthographic inconsistency, and morphological complexity. Findings from this study affirm that transfer learning from related languages offers practical advantages but cannot fully substitute for in-domain, dialectally representative datasets. This study also suggests the need for morphologically aware and dialect-sensitive modeling approaches to improve ASR accuracy for languages such as Ewe. Future work should prioritize leveraging this study's unlabeled speech corpus through self-supervised learning techniques and explore domain-adapted language modeling to enhance transcription reliability for critical applications such as healthcare, education, and public service delivery. By addressing these linguistic and technological gaps, this research lays the foundation for more inclusive speech technologies that preserve and promote the use of indigenous African languages. The data splits and trained model is publicly available on GitHub and Huggingface.

## Limitations

While self-supervised learning approaches, such as wav2vec 2.0 (Baevski et al., 2020), have shown promise in reducing the dependence on labeled data by leveraging large amounts of unlabeled audio, their application is not without challenges. Although the dataset collected in this study comprises 900 hours of unlabeled Ewe speech, the computational constraints limited the feasibility of training with wav2vec. Access to large-scale computing resources remains a significant bottleneck in LRE research. This study argues that, in LREs, model performance is fundamentally constrained by linguistic complexity and computational resources rather than modeling innovations alone. Addressing these challenges is essential to advancing equitable access to speech technologies for underrepresented languages. Future research may build on these findings by prioritizing the development of scalable methodologies and resources that enable the advancement of ASR technologies for LRLs, such as Ewe and other Ghanaian languages.

## Ethics Statement

Ethical approval for this study was obtained from the Ethics Committee for Basic and Applied Sciences, University of Ghana. All participants, including recorders, validators, and transcribers, were informed of the study's goal and that the collected data would be used for research purposes only. Also, the consent form specified that participation was voluntary and participants could withdraw at any point. All participants were compensated for their respective contributions.

## Acknowledgement

# References

Kwame Badu Antwi-Boasiako and Kofi Agyekum. 2022. Globalization, colonization, and linguicide: How ghana is losing its local languages through radio and television broadcast. *International Journal of Humanities and Social Science*, 12(8):142–151.

Rosana Ardila, Megan Branson, Kelly Davis, Michael Kohler, Josh Meyer, Michael Henretty, Reuben Morais, Lindsay Saunders, Francis Tyers, and Gregor Weber. 2020. Common voice: A massively-multilingual speech corpus. In *Proceedings of the Twelfth Language Resources and Evaluation Conference*, pages 4218–4222, Marseille, France. European Language Resources Association.

Alexei Baevski, Henry Zhou, Abdelrahman Mohamed, and Michael Auli. 2020. wav2vec 2.0: A framework for self-supervised learning of speech representations. In *Proceedings of the 34th Conference on Neural Information Processing Systems (NeurIPS 2020)*, pages 12449–12460.

Laurent Besacier, Etienne Barnard, Alexey Karpov, and Tanja Schultz. 2014. Automatic speech recognition for under-resourced languages: A survey. *Speech Communication*, 56(1):85–100.

Lindell Bromham, Russell Dinnage, Hedvig Skirgård, Andrew Ritchie, Marcel Cardillo, Felicity Meakins, Simon Greenhill, and Xia Hua. 2022. Global predictors of language endangerment and the future of linguistic diversity. *Nature Ecology & Evolution*, 6(2):163–173.

Chris Callison-Burch and Mark Dredze. 2010. Creating speech and language data with amazon's mechanical turk. In *Proceedings of the Workshop on Creating Speech and Language Data with Amazon's Mechanical Turk*, pages 1–12.

Ilaria Chizzoni and Alessandro Vietti. 2024. Towards an ASR system for documenting endangered languages: A preliminary study on sardinian. Technical report, Free University of Bozen-Bolzano.

De Graft Johnson Dei. 2024. Sustainability and development of Ewe communities in ghana through indigenous knowledge management practices. *Collection and Curation*, 43(4):111–123.

Kavan Fatehi, Mercedes Torres Torres, and Ayse Kucukyilmaz. 2025. An overview of high-resource automatic speech recognition methods and their empirical evaluation in low-resource environments. *Speech Communication*, 167:103151.

Alexandru-Lucian Georgescu, Horia Cucu, Andi Buzo, and Corneliu Burileanu. 2020. RSC: A Romanian read speech corpus for automatic speech recognition. In *Proceedings of the Twelfth Language Resources and Evaluation Conference*, pages 6606–6612, Marseille, France. European Language Resources Association.

Alexander Gutkin, Isin Demirsahin, Oddur Kjartansson, Clara Rivera, and Kólá Túbòsún. 2020. Developing an open-source corpus of yoruba speech. In *Proceedings of Interspeech 2020*, pages 404–408.

Umar Adam Ibrahim, Moussa Mahamat Boukar, and Muhammed Aliyu Suleiman. 2022. Development of Hausa dataset: A baseline for speech recognition. *Data in Brief*, 40:107820.

Robbie Jimerson, Kruthika Simha, Raymond Ptucha, and Emily Prud'hommeaux. 2018. Improving ASR output for endangered language documentation. In *Proceedings of the 6th Workshop on Spoken Language Technologies for Under-Resourced Languages (SLTU 2018)*, pages 187–191.

Minu Kim, Kangwook Jang, and Hoirin Kim. 2025. Improving cross-lingual phonetic representation of low-resource languages through language similarity analysis. arXiv preprint arXiv:2502.01234.

Mark Atta Mensah, Isaac Wiafe, Akon Ekpezu, Kwame Appati, Jamal-Deen Abdulai, Akosua Nyarkoa Wiafe-Akenten, Frank Ernest Yeboah, and Gifty Odame. 2025. Benchmarking akan asr models across domain-specific datasets: A comparative evaluation of performance, scalability, and adaptability. In *Accepted - Future Technologies Conference*.

Tobi Olatunji, Tejumade Afonja, Aditya Yadavalli, Chris Chinenye Emezue, Sahib Singh, Bonaventure Francis-Pierre Dossou, Joanne Osuchukwu, Salomey Osei, Atnafu Lambebo Tonja, Naome Etori, and Clinton Mbataku. 2023. AfriSpeech-200: Pan-african accented speech dataset for clinical and general domain ASR. *Transactions of the Association for Computational Linguistics*, 11:1669–1685.

Vassil Panayotov, Guoguo Chen, Daniel Povey, and Sanjeev Khudanpur. 2015. Librispeech: An ASR corpus based on public-domain audio books. In *Proceedings of the 2015 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP 2015)*, pages 5206–5210.

Kay Peterson, Audrey Tong, and Yan Yu. 2021. OpenASR20: An open challenge for automatic speech recognition of conversational telephone speech in low-resource languages. In *Proceedings of Interspeech 2021*, pages 4324–4328.

Alec Radford, Jong Wook Kim, Tao Xu, Greg Brockman, Christine McLeavey, and Ilya Sutskever. 2022. Robust speech recognition via

large-scale weak supervision. In *Proceedings of the 40th International Conference on Machine Learning (ICML '23)*, pages 28492–28518.

Alessandro Ragano, Emmanouil Benetos, and Andrew Hines. 2020. Development of a speech quality database under uncontrolled conditions. In *Proceedings of Interspeech 2020*, pages 4616–4620.

Philip Resnik, Mari Broman Olsen, and Mona Diab. 1999. The bible as a parallel corpus: Annotating the 'book of 2000 tongues'. *Computers and the Humanities*, 33(1):129–153.

Jemima Sam and Cynthia Ablah Agbloe. 2024. Lexical variations in the Ewe language spoken in ho in the volta region of ghana. *Journal of Education*, (7):69–87.

Tanja Schultz. 2002. Globalphone: A multilingual speech and text database developed at karlsruhe university. In *Proceedings of the 7th International Conference on Spoken Language Processing (ICSLP 2002)*, pages 345–348.

Mohammad Teduh Uliniansyah, Gunarso, Elvira Nurfadhilah, Lyla Ruslana Aini, Juliati Junde, Fara Ayuningtyas, and Agung Santosa. 2016. A tool to solve sentence segmentation problems on preparing a speech database for an indonesian text-to-speech system. In *Procedia Computer Science*, volume 81, pages 188–193.

Isaac Wiafe, Jamal-Deen Abdulai, Akon Obu Ekpezu, Raynard Dodzi Helegah, Elikem Doe Atsakpo, Charles Nutrokpor, Fiifi Baffoe Payin Winful, and Kafui Kwashie Solaga. 2025. Advancing automatic speech recognition for low-resource ghanaian languages: Audio datasets for akan, ewe, dagbani, dagaare, and ikposo. *Data in Brief*, 61:111880.

# CLEAR: Code-Mixed ASR with LLM-Driven Rescoring

**Shivam Kumar, Md. Shad Akhtar**
Department of CSE, IIIT-Delhi, New Delhi, India
{shivam23004,shad.akhtar}@iiitd.ac.in

## Abstract

Code-mixing presents significant challenges for Automatic Speech Recognition (ASR), especially for Indian languages, due to homophone ambiguity, domain-specific word identification, and data scarcity. Traditional ASR models struggle with these complexities, often failing to differentiate between phonetically similar words in multilingual contexts. To address this, we propose CLEAR, a novel rescoring model that integrates descriptive prompting and LLM-based rescoring while analyzing the impact of $n$-best hypotheses across multiple beam widths. CLEAR enhances ASR performance, achieving S-WER of 26.9, P-WER of 26.46, and T-WER of 25.04—improving by 6.9%, 13.47%, and 4.42%, respectively, over the best baseline TDNN. These findings demonstrate that CLEAR effectively resolves homophone ambiguities and refines transcriptions, leading to a 13.56% S-WER and 7.77% T-WER reduction over decoder only fine-tuned Whisper.

## 1 Introduction

Code-mixing and code-switching are prevalent linguistic phenomenon in multilingual communities, where speakers alternate between languages within a single discourse. The terms code-switching and code-mixing are often used interchangeably[1]; however, they carry distinct linguistic meanings. Code-switching is inter-sentential where the language/code switch occurs at the utterance level; while code-mixing is intrasentential where the switch happens at the word or phrase level within a sentence (Thara and Poornachandran, 2018; Setiawan, 2023; Winata et al., 2022).

Developing automatic speech recognition (ASR) systems for code-switched speech



Figure 1: Proposed architecture of CLEAR.

presents unique challenges due to the nature of the language(Çetinoğlu et al., 2016). Code-switching introduces linguistic complexities such as cross-lingual homophone disambiguation (e.g., Bill in English means a receipt or a piece of paper; however, बिल(Bill) in Hindi means a hole or a burrow) (Yu et al., 2024), code-switching point detection (Wang et al., 2019), and the identification of embedding and matrix languages, which are crucial for determining the correct syntactic structure (e.g., Subject-Verb-Object (SVO) in English vs. Subject-Object-Verb (SOV) in Hindi) (Iakovenko and Hain, 2024). Accurately modeling these aspects is essential for generating grammatically coherent transcriptions.

To mitigate these challenges, recent advancements (Prabhavalkar et al., 2023) in ASR have been driven by innovations in neural architectures, training strategies, and robust learning techniques. Many explorations in end-to-end ASR leveraged transformer-based encoders and self-supervised pretraining to learn rich representations from raw acoustic signals (Baevski et al., 2020; Hsu et al., 2021; Chadha et al., 2022). Whisper (Radford

---

[1]We will also use these two terms interchangeably to refer our input setting in this paper.

et al., 2023) is one of most popular state-of-the-art multilingual ASR models, trained on 630K hours of data, that can transcribe, translate, and detect speeches across 99 languages. Peng et al., 2023 incorporated a mixture of language tags, i.e., `<zh><en>` for Mandarin-English code-switched dataset; while Yang et al., 2024c extended Whisper with a separate language tag `<|en−zh|>` to uniquely refer to Mandarin-English code-mixed sentences.

With the inception of LLMs, recent studies have leveraged prompting strategies (Liu et al., 2023b) in Whisper-based architectures. Suh et al., 2024 explored LLMs to generate contextual descriptions, which were then used to prompt the Whisper for transcription, demonstrating the potential of LLM-guided ASR.

Although previous methods produced good results, they required generating different mixes of language tags or using an LLM to create prompts for each utterance, we do not aim to do that also these works have been experimented on Mandarin-English code-switchted dataset belongs to different domain and not on Hindi-English code-switched which has its own linguistic complexities. Our research builds upon these works but takes a distinct approach. Unlike prior studies, we do not introduce new language tags or rely on LLM-generated prompts. Instead, we propose CLEAR to demonstrate that descriptive prompting alone can yield high-quality code-mixed transcriptions by refining ASR outputs through LLM-based rescoring. We hypothesize that LLM-based rescoring can mitigate the issue of homophone disambiguation by scoring the fluent sentences with higher scores. We obtain $n$-beams of potential outputs from Whisper and utilize LLM-based scorer to measure their linguistic fluency and coherency. We fine-tune the Whisper decoder while keeping the encoder in a frozen state. Specifically, we experiment with leading LLMs – GPT-2, LLaMA 3.1 (8B), LLaMA 3.2 (1B), DeepSeek R1[2], Qwen-2 (7B), Mistral (7B), and GPT-4 (Radford et al., 2019; Dubey et al., 2024; Guo et al., 2025; Yang et al., 2024a; Jiang et al., 2023; Achiam et al., 2023) – while varying beam widths to assess their impact on ASR performance. We employ MUCS 2021 dataset

---

[2]DeepSeek-R1-Distill-Llama-8B

to evaluate CLEAR. Our analysis across six competitive baselines signifies the importance of CLEAR on code-switched ASR outputs. Our contributions are summarized below:

- We present CLEAR and contribute novel insights into the role of descriptive prompting and LLM-based scoring in improving code-switched ASR systems, paving the way for more effective transcription models in multilingual settings.

- We extensively evaluate CLEAR against 6 baseline methods. We present CLEAR significant improvement evaluated across evaluation metrics.

**Reproducibility:** https://github.com/flamenlp/CLEAR

## 2 Related Work

Another line of work have explored approaches to improve language modeling and context understanding. For instance, Aditya et al., 2024 investigated attention mechanisms within transformer layers, identifying attention heads that effectively capture language identities and guiding them accordingly. To mitigate multilingual context confusion, Zhang et al., 2022 proposed attention weight recomputation to better differentiate languages within speech. Further, Liu et al., 2023a, 2024a introduced language biases at both the token and frame levels to enhance the model's ability to handle language switching effectively. Song et al., 2022 proposed a language-specific characteristic assistance (LSCA) method to mitigate the problem caused by lanuage-specific encoders (LSEs) since most existing methods did not have language constraints; they introduced a language-specific loss to do that. To disambiguate homophones Srivastava and Sitaram, 2018 used a WX-based common pronunciation scheme for mixed language pairs and unification of homophones during training, resulting in a lower word error rate for systems built using this data. Chung et al., 2022 proposed a novel homophone extension method to integrate human knowledge of the homophone lexicon into the beam search decoding process with language model re-scoring. Some of the recent work has also explored Mixture of Experts

(MoE) architectures for code-switched ASR. Ye et al., 2024 proposed using separate encoders as language experts, while Yang et al., 2024b introduced a disentanglement loss to enable lower encoder layers to capture interlingual acoustic information while reducing linguistic confusion in higher layers. Liu et al., 2024b introduced a language alignment loss in ASR training to align acoustic features to pseudo-language labels learned from the decoder and also employs LLM via generative error correction to tackle the problem caused by complex token alternatives for language modeling in bilingual scenarios.

## 3 Methodology

In this section, we describe **CLEAR** model to enhance transcription of Hindi-English code-switched speech. Our primary objective is to use a descriptive prompting strategy that provides contextual guidance to the decoder, improving transcription accuracy without requiring extensive fine-tuning of the entire model. Additionally, we incorporate LLMs for rescoring to further refine the final transcription output.

**Proposed Pipeline:** Whisper, unlike conventional ASR models, allows for prompting (Suh et al., 2024) through special tokens that guide the transcription process. These tokens include `<|sop|>` (*start-of-previous*), `<|sot|>` (*start-of-transcript*), `<|en|>` or `<|hi|>` (language tags), `<|transcribe|>` (specifies the task as transcription), and `<|notimestamps|>` (to disable word-level timestamps). In our proposed pipeline, we strategically place our custom prompt after the `<|sop|>` token. This placement provides contextual information to the decoder while ensuring compliance with Whisper's input constraints, as only 224 tokens[3] can be used as a prompt. We also experiment with different prompts to check which is working best, more details will be discussed in section 4. Our constructed prompt follows the format:

`<|sop|><|prompt|><|hi|><|transcribe|><|notimestamps|>`

We fine-tune Whisper by integrating our designed prompt to influence the decoder's

---

behavior.

**LLM-Based Scorer:** To further refine transcription quality, we introduce a rescoring mechanism utilizing LLMs. LLMs are trained on vast multilingual corpora (Gurgurov et al., 2024), effectively capture semantic structures and can assist in selecting the most plausible transcription candidate. Given a beam width of $n$, the Whisper decoder generates $n$ transcription hypotheses $\{x^{(1)}, x^{(2)}, \ldots, x^{(n)}\}$. The rescoring process involves computing the sum of log probability of each hypothesis $x^{(i)}$ based on the LLMs ouput logits, we call this sum as score. The log probability of a candidate sequence is given by:

$$
\begin{aligned}
\log P(\mathbf{x}^{(i)}) &= \sum_{t=1}^{T} \log P(x_t | x_{1:t-1}, \theta) \\
&= \sum_{t=1}^{T} \log \left( \frac{\exp(z_t^{(x_t)})}{\sum_j \exp(z_t^{(j)})} \right) \\
&= \sum_{t=1}^{T} \left( \log \left( \exp(z_t^{(x_t)}) \right) - \log \sum_j \exp(z_t^{(j)}) \right) \\
&= \sum_{t=1}^{T} \left( z_t^{(x_t)} - \log \sum_j \exp(z_t^{(j)}) \right)
\end{aligned}
\tag{1}
$$

where $t$ is current time step, and $T$ is the total number of steps. The sequence $x = (x_1, x_2, \ldots, x_T)$ consists of tokens $x_t$, each assigned a logit $z_t^{(x_t)}$ by the LLM. The LLMs parameters are denoted by $\theta$. $j$ is the index in the output vocabulary.

The best transcription $x^*$ is then selected as the one with the highest log probability among all n hypotheses:

$$
\mathbf{x}^* = \arg \max_{i \in 1,2,\ldots,n} \log P(\mathbf{x}^{(i)})
\tag{2}
$$

We conduct experiments with different LLMs and beam widths to assess their impact on transcription quality.

**Fine-Tuning:** Our fine-tuning process focuses solely on the decoder while keeping the encoder frozen. Since Whisper has been pretrained on 630K hours of multilingual speech data, its encoder already possesses a strong understanding of the acoustic properties of speech. Freezing the encoder prevents overfitting and ensures that the model retains its

| Prompt-1 | Prompt-2 | Prompt-3 |
|---|---|---|
| This transcript is a code-switched text. Mix of devnagri and english words are present. Text is related to tutorials on academic or technical subjects. | This transcript is a code-switched text. Mix of Devanagari and English words are present. Text is related to tutorials on academic or technical subjects. Few examples look like this: तो electricity bill option पर click करें कुछ गलत password दीजिए और enter प्रेस करें | The transcript comprises telephone quality speech data in Hindi. Transcript is mixed of Hindi and English words like this: hi hi hi hi en hi hi hi hi hi en hi hi hi hi. Transcribe the speech in this format. |

Table 1: Prompts use to fine-tune the CLEAR is listed here. We experiment with many prompts, few of them are shown in this table.

general ASR capabilities while adapting its decoder for improved handling of code-switched speech. Previous studies (Yang et al., 2023) have shown that Whisper with a frozen encoder can achieve superior performance on certain ASR tasks. We will also prove that decoder-only finetuning works in the section 5. Our methodology, which integrates Whisper's ASR capabilities with descriptive prompting and LLM-based rescoring, presents an efficient approach for improving code-switched speech recognition without altering language tags (Yang et al., 2024c).

## 4 Experiments

**Dataset Description:** We use the Hindi-English code-switched dataset[4] from the MUCS Challenge, Interspeech 2021 (Diwan et al., 2021), derived from spoken tutorials on technical topics. All audio files are sampled at 16 kHz with 16-bit encoding. The dataset comprises spontaneous speech from educational settings, making it particularly challenging due to variations in speaker accents, speech disfluencies, and technical terms. It consists of ~100 hours of data and splitted into train (89.86 hrs), test (5.18 hrs), and blind (6.24 hrs) sets. Moreover, the dataset has code-switching percentage[5] of 85.88%, 81.88%, and 95.55% in train, test, and blind sets, respectively. The dataset also contains enunciated punctuation (e.g., "<"

for "lesser"). Furthermore, the train and test sets have ~33.9% overlaps; however, blind test-train sentence overlap is 2.1%. Therefore, we evaluate CLEAR on blind set only.

**Evaluation Metric:** For evaluation, we compute three variants of word-error-rate (WER), i.e., strict-WER (S-WER), punctuation-WER (P-WER), and transliterated-WER (T-WER). S-WER is the standard WER metric which computes the transcription error rate at the surface level. In comparison, P-WER accounts for variations in punctuation by applying a predefined punctuation mapping before computing WER (replace "greater" by ">"); thus, ensuring consistency in evaluation. On the other hand, T-WER assesses errors in cross-linguistic transcription by replacing words in the predicted transcriptions with their corresponding transliterated forms (replace डिस्क्रिप्शन by description). These three WER variants provide a comprehensive assessment of the ASR model, capturing both standard transcription errors and linguistic variations of code-mixed languages.

**Descriptive Prompt Details:** We experiment with different numbers of prompts in different styles to guide the decoder; a few of them worked and some did not. Some examples of prompts are listed in Table 1. The prompts that gives the best performance are Prompt-1 and Prompt-2 and also in in CLEAR architecure we have used this Prompt-1, you can see in the Fig. 1. Prompt-1 is very simple

---

[4]https://www.openslr.org/104/
[5]code-switched utterances upon total no. of utterances

| Model | S-WER | P-WER | T-WER |
|---|---|---|---|
| TDNN | 28.90 | – | 26.20 |
| E2E Transformer | 33.65 | – | 29.80 |
| PromptingWhisper | 42.10 | 41.72 | 34.75 |
| Whisper (ZS) | 266.5 | 265.84 | 263.30 |
| Whisper (FT) | 32.16 | 31.64 | 28.54 |
| Whisper (FE) | 31.12 | 30.58 | 27.15 |
| CLEAR | **26.9** (↓ 6.9%) | **26.46** (↓ 13.47%) | **25.04** (↓ 4.42%) |

Table 2: Comparative results on MUCS dataset. ZS = Zero Shot, FT = Full finetuning, and FE = Finetune by Frozen Encoder

to tell the decoder that text is code-switched and guiding what code-switching means by mentioning about type of languages involed and also telling the decoder that dataset is related to technical subject and tutorial taught in academic setting, basically it is giving the contextual cues about the dataset. Similary Prompt-2 is doing the same, additionally it is also giving the examples to make it more clear. And Prompt-3 even tells the decoder when the language switch is happening, which will allow the decoder to predict the language switch even better, one of the major problems of code-switching. But sometimes these Prompt-2, 3 and some more which are not listed here do not work. One of the problems we find is that these are longer prompts, which may be longer for some utterances, and Whisper can understand only 224 tokens as a prompt and anything longer that will be truncated. Therefore, Prompt-3 mostly will fail due to its dynamic nature for each utterance, Prompt-2 can work if it is within the token limits.

**Baselines:** We compare the performance of our CLEAR model with following baselines. DNN-HMM (Diwan et al., 2021) is a neural network created using the Kaldi toolkit[6], consisting of 8 TDNN (Time-Delay Neural Network) blocks (Peddinti et al., 2015) with a dimension of 768. End-to-End (E2E) Transformer (Diwan et al., 2021) is a hybrid CTC-Attention model (Watanabe et al., 2017) with a 12-layer encoder and a 6-layer decoder, each with 2048 units and 0.1 dropout rate. It employs a CTC weight of 0.3 and an attention weight of 0.7, using eight 64-dimensional

---

| Scorer | Beam | S-WER | P-WER | T-WER |
|---|---|---|---|---|
| – | 1 | 28.09 (↓ 9.73%) | 28.42 (↓ 7.06%) | 26.06 (↓ 4.01%) |
| GPT2 | 5* | **26.9** | **26.46** | **25.04** |
|  | 10 | 27.11 | 26.67 | 25.47 |
| LLaMA 3.1 (8B) | 5 | 28.26 | 27.78 | 25.02 |
|  | 10 | 28.39 | 27.91 | 25.09 |
| LLaMA 3.2 (1B) | 5 | 27.75 | 27.30 | 25.69 |
|  | 10 | 27.86 | 27.42 | 26.03 |
| DeepSeek | 5 | 27.60 | 27.14 | 25.55 |
|  | 10 | 28.20 | 27.74 | 26.35 |
| Qwen-2 (7B) | 5 | 27.48 | 27.01 | 25.40 |
|  | 10 | 27.89 | 27.44 | 26.06 |
| Mistral (7B) | 5 | 27.55 | 27.11 | 25.49 |
|  | 10 | 27.91 | 27.50 | 26.09 |
| GPT-4 | 5 | 27.85 | 27.29 | 24.82 |

Table 3: Ablation on different beam widths. Star(*) signifies the CLEAR model.

attention heads per layer. PromptingWhisper (Peng et al., 2023) is a Whisper-large model tested in a zero-shot setup by changing the language tag. Here, we use `<|hi|><|en|>` as the language tag. Additionally, we employ Whisper in both zero-shot and fine-tuning settings. For CLEAR we utilize Whisper-small[7] (12 self-attention layers in both the encoder and the decoder) which is capable of processing 30-second audio segments and generates text autoregressively.

**Training Details:** We fine-tune CLEAR for 10 epochs using a learning rate of 1e−4. The training was conducted on an NVIDIA A100 GPU with a batch size of 16. We employ AdamW as the optimizer with a weight decay of 0.01 to regulate parameter updates. The training procedure leverages mixed-precision training to improve efficiency and reduce memory consumption while maintaining numerical stability. During inference, we adopt beam search decoding (beam-width $= n$) to improve transcription accuracy. We also experiment with temperature scaling to prevent overly confident incorrect predictions.

## 5  Results and Analysis

Table 2 presents the comparative analysis of CLEAR against other baselines on the blind set of the MUCS dataset. Our initial evaluation

---

of Whisper in a zero-shot (ZS) setting revealed a drastic degradation in performance, with S-WER exceeding 260%. This excessively high WER is primarily caused by repetitive character sequences and a lack of domain adaptation, leading to significant transcription errors. The model struggles with both code-switching regions and the specific linguistic patterns of the dataset, underscoring the limitations of the out-of-the-box Whisper-small model in handling code-mixed speech. To mitigate these issues, we fine-tune Whisper in two settings, first full fine-tuning (FT) and second fine-tuning by frozing encoder (FE), results are shown in Table 2. Whisper (FT) shows improvement in S-WER (32.16), P-WER (31.64), and T-WER (28.54) as compared to ZS setting. We further fine-tune Whisper by frozing encoder (FE) on our dataset, allowing it to adapt to the linguistic characteristics of code-switched speech. This results in a substantial improvement in S-WER (31.12) compared to ZS setting and this is also an improvement over Whisper(FT), as the model exhibit a better understanding of language semantics, reduced character repetition, and improved handling of code-switching boundaries. The results of Whisper(FT) and Whipser(FE) shows that fine-tuning by the frozing encoder works very well as compared to full fine-tuning. Therefore, for our proposed pipeline we will use Whipser (FE) for all our experiments. Despite these gains, domain-specific challenges persisted, particularly for low-frequency words, homophone inconsistency, and unseen terms that were not well represented in the corpus.

To further enhance transcription accuracy, we use descriptive prompt-based fine-tuning approach (Peng et al., 2023), where the Whisper decoder was fine-tuned while keeping the encoder frozen. This strategy led to a 9.73% reduction in S-WER and a 4.01% reduction in T-WER (beam $n = 1$ in Table 3) compared to Whisper (FE). These results indicate that prompt-based fine-tuning is an effective adaptation strategy for code-switching ASR, significantly improving performance without requiring architectural modifications to the Whisper model. Building on this, we employ our scoring mechanism to grade the $n$-best hypothesis from the fine-tuned Whisper model. CLEAR

reports a reduced score of S-WER (26.9), P-WER (26.46), and T-WER (25.04) – a significant reduction of +6.9%, +13.47% and +4.42% in S-WER, P-WER, and T-WER, respectively, over the best baseline (i.e TDNN).

As reported in Table 3, we employ multiple open-source LLMs, such as GPT-2, LLaMA 3.1 (8B), LLaMA 3.2 (1B), DeepSeek, Qwen-2 (7B), Mistral (7B), and GPT-4 for transcription scoring. We systematically analyzed the impact of beam width on WER, experimenting with $n = 5, 10, 15$ and $20$[8]. Our findings indicate that a beam width of 5 consistently yielded the best results across models. In particular, compared to fine-tune (FE) Whisper-small, GPT-2 achieved a 13.56% reduction in strict WER and a 7.77% reduction in transliterated WER; LLaMA 3.1 (8B) yielded a 9.19% reduction in strict WER and a 7.84% reduction in transliterated WER; LLaMA 3.2 (1B) yielded a 10.82% reduction in strict WER and a 5.37% reduction in transliterated WER; Deepseek yielded a 11.31% reduction in strict WER and a 5.89% reduction in transliterated WER; Qwen yielded a 11.69% reduction in strict WER and a 6.44% reduction in transliterated WER; Mistral yielded a 11.47% reduction in strict WER and a 6.11% reduction in transliterated WER; GPT-4 yielded a 10.50% reduction in strict WER and a 20.24% reduction in transliterated WER, It outperforms GPT2 by 0.87% only in T-WER. However, we would like to highlight that GPT-4 incurred a significant 3\$/100 hypotheses during inference. This justifies the use of GPT-2 based scorer in CLEAR as a budget effective solution. Also rescoring mechanism with GPT-4 was different compared to other LLMs because we can't extract the logits from GPT-4 therefore we rank the hypothesis based on their accuracy, coherency, and fluency on a scale of $-10$ (very inaccurate) to 10 (perfectly accurate). Calculating sum of log probabilites of logits will not work in the case of GPT-4.
These results highlight the effectiveness of LLM-based rescoring techniques, demonstrating their ability to refine transcription outputs and further reduce errors in code-mixed ASR

---

[8]We do not report n=15 and n=20 due to inferior results. We observe performance drops across multiple scorers with beam>5.

| | Model | Text | Remarks |
|---|---|---|---|
| **Sentence 1** | GT | अब इस method पर आते हैं | – |
| | W-(ZS) | अब इस मेफ पर आते | "method" is missing, and substitution leads to loss of meaning |
| | W-(FE) | और formula का स्टेटमेंट दें | Extraneous words introduced, altering the intended meaning. |
| | CLEAR–R | अब इस method पर आते हैं | Matches GT exactly, correct transcription. |
| | CLEAR | अब इस method पर आते हैं | Matches GT exactly, correct transcription. |
| **Sentence 2** | GT | get noise profile पर click करें | – |
| | W-(ZS) | अगर तो तो तो तो तो तो तो तो | Repeated words result in meaningless output. |
| | W-(FE) | cat noise profile पर click करें | Incorrect word "cat" instead of "get", missing domain-specific knowledge. |
| | CLEAR–R | get noise profile पर click करें | Correctly retains "get" and follows GT structure. |
| | CLEAR | get noise profile पर click करें | Matches GT exactly, preserving domain-specific knowledge. |
| **Sentence 3** | GT | अब वापस IDE पर आते हैं | – |
| | W-(ZS) | अब वापस आईडी पर आते | Misrecognition of "IDE" as **"आईडी"** changes the meaning. |
| | W-(FE) | अब वापस ID पर आते हैं | Homophones is not correctly identified |
| | CLEAR–R | अब वापस IDE पर आते हैं | Correctly identifies "IDE" and follows GT. |
| | CLEAR | अब वापस IDE पर आते हैं | Matches GT exactly, resolving homophone ambiguity |
| **Sentence 4** | GT | चिंता न करें यदि class diagram view में नहीं खुलता है | – |
| | W-(ZS) | जितना न करें ये दिखाएगा जगे व्यू में नहीं खुलता है | Unintelligible phrase with extra words. |
| | W-(FE) | चिंता न करें यदि, टगग्राम डायग्राम में नहीं खुलता है | Incorrectly replaces "class diagram view" with an unrelated term. |
| | CLEAR–R | चिंता न करें यदि class tag view में नहीं खुलता है | Partially correct, but "class tag view" is incorrect. |
| | CLEAR | चिंता न करें यदि class diagram view में नहीं खुलता है | CLEAR Matches GT exactly, ensuring correct code-switching. |
| **Sentence 5** | GT | 1123 put insulin . fasta file के लिए contents | – |
| | W-(ZS) | अप्रोट इंसुलिन ड़ प्रश्टा फाँईल के लिए, खन्टेंस दिखाता है. | Misrecognized words distort the sentence meaning. |
| | W-(FE) | आउटपुट insulin dot first फाइल के लिए कंटेंट्स दिखाता है | Incorrect segmentation of "fasta file" as "dot first file." |
| | CLEAR –R | default insulin dot firster file के लिए कंटेंट्स दिखाता है | Better than W-(FE), but still modifies "fasta file." |
| | CLEAR | default installation dot firster file के लिए कंटेंट्स दिखाता है | CLEAR Corrects some errors but still alters "fasta file." |

Table 4: Comparison of ASR outputs among competitive models. CLEAR–R: CLEAR without scorer (i.e., beam=1 and Whisper fine-tuned with descriptive prompt).

tasks. Overall, our findings establish that a combination of prompt-based fine-tuning and LLM-based rescoring substantially enhances the performance of the code-mixed ASR task.

**Complexity of LLMs:** While large language models (LLMs) are often associated with high computational costs, we carefully designed CLEAR to avoid these burdens. Rather than fine-tuning the LLMs, which would require significant GPU hours and memory, we utilized them solely for inference to score ASR hypotheses. This approach is lightweight - each hypothesis takes only about 0.6 to 0.7 seconds to evaluate — making it both efficient and practical for real-world post-processing scenarios. Interestingly, we also observed significant improvements in transcription quality even when using relatively smaller LLMs such as LLaMA-3.2 (1B). This suggests that even modest-sized language models can capture contextual nuances well enough to resolve ambiguities and correct ASR output, particularly in challenging settings such as code-mixed speech. This balance between performance gain and computational overhead

suggests that we can use LLMs as scorers for post-processing the ASR outputs and makes our method feasible for broader deployment.

**Qualitative and Error Analysis:**
To further assess the quality of the generated transcriptions, we conduct a detailed qualitative and error analysis, as shown in Table 4. This analysis highlights the advantages of our proposed pipeline over the zero-shot and fine-tuned models in handling various linguistic complexities in code-switched ASR. Our pipeline exhibits significant improvements in multiple aspects, including word ordering, domain-specific terminology, homophone disambiguation, and rescoring-based refinement.

- **Word Ordering:** One notable enhancement is the model's ability to predict words in the correct order. As observed in Sentence 1, while the zero-shot model produces an incomplete and incorrect phrase, and the fine-tuned model introduces extraneous words, both the CLEAR–R and CLEAR versions successfully reconstruct the correct sentence structure. This sug-

gests that our approach effectively learns the sequential dependencies within code-mixed speech, leading to more coherent and grammatically accurate outputs.

- **Domain-specific terminology:** Another key improvement is in handling domain-specific terminology, as illustrated in Sentence 2. The fine-tuned model incorrectly transcribes "cat" instead of "get", demonstrating its struggle with differentiating both similar-sounding words and domain-specific words. CLEAR–R approach, however, accurately transcribes "get", showcasing a better understanding of contextual cues. This improvement is crucial in technical and specialized domains, where precise word recognition significantly impacts the usability of transcriptions.

- **Homophone Disambiguation:** In Sentence 3, CLEAR effectively tackles homophone disambiguation, particularly distinguishing between "id" and "ide". The fine-tuned model fails to capture this distinction, incorrectly predicting "id" instead of "ide", whereas the CLEAR–R and CLEAR accurately recognize the correct term. This capability is essential in technical environments where similar-sounding words carry distinct meanings, ensuring that transcripts remain contextually relevant.

- **Rescoring Refinements:** In sentence 4, demonstrates the effectiveness of rescoring in refining transcriptions. While the CLEAR–R output closely aligns with the ground truth, it incorrectly transcribes "class tag view" instead of "class diagram view". The CLEAR output successfully corrects this error, emphasizing the impact of our rescoring mechanism in enhancing transcription accuracy. This step ensures that even when the initial transcription is suboptimal, the model can refine its predictions to achieve greater alignment.

- **Pitfall of CLEAR:** It is not like our CLEAR model is giving good performance of every utterances. In some scenarios it is failing. As we can see in the sentence 5, it is not correctly predicting many words. One

reason for this may there are two many code-switches available, puctuation "." is present, and domain specific word is also present. There may be possibility because the presence of these challenges causing the problem to the ASR.

## 6 Conclusion

This study introduces CLEAR, a novel approach to enhancing Hindi-English code-mixed ASR by integrating descriptive prompt-based fine-tuning and LLM-based rescoring. Our findings reveal that fine-tuning the Whisper decoder while freezing the encoder is a highly effective strategy for code-switching transcription (refer $2^{nd}$ last row of Table 2), yielding substantial reductions in various word error rates.Extensive ablation and qualitative analysis establishes LLM-based rescoring as an efficient refinement mechanism, effectively disambiguating homophones, and also enhances overall readability and domain-specific accuracy without requiring explicit language tags or specialized pretraining. These insights pave the way for more adaptable and resource-efficient LLM-guided ASR systems, particularly in low-resource and multilingual settings.

## Limitations

This work has also certain limitations. Due to GPU memory constraints, we did not explore the impact of larger batch sizes, which could potentially influence model performance. Additionally, there is a lack of a diverse and high-quality code-mixed dataset for Indian languages. Our experiments are limited to MUCS code-switching data. To the best of our knowledge, this is the only publicly available dataset. We did not evaluate the approach on other pairs of Indian languages. The limited dataset diversity hinders the robustness, highlighting the critical need for high-quality datasets in this research area. However, we believe that the overall pipeline of our CLEAR model is highly adaptable to other language pairs as well. Another challenge is tackling too many code-switching points, availability of punctuation, and domain specific words all these togther if present in a sentence. sometimes CLEAR fails to handle this. In future we also plan to handle this limitaion of our model.

## Acknowledgment

## References

Josh Achiam, Steven Adler, Sandhini Agarwal, Lama Ahmad, Ilge Akkaya, Florencia Leoni Aleman, Diogo Almeida, Janko Altenschmidt, Sam Altman, Shyamal Anadkat, et al. 2023. Gpt-4 technical report. *arXiv preprint arXiv:2303.08774*.

Bobbi Aditya, Mahdin Rohmatillah, Liang-Hsuan Tai, and Jen-Tzung Chien. 2024. Attention-guided adaptation for code-switching speech recognition. In *ICASSP 2024-2024 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, pages 10256–10260. IEEE.

Alexei Baevski, Henry Zhou, Abdelrahman Mohamed, and Michael Auli. 2020. wav2vec 2.0: A framework for self-supervised learning of speech representations.

Özlem Çetinoğlu, Sarah Schulz, and Ngoc Thang Vu. 2016. Challenges of computational processing of code-switching. *arXiv preprint arXiv:1610.02213*.

Harveen Singh Chadha, Priyanshi Shah, Ankur Dhuriya, Neeraj Chhimwal, Anirudh Gupta, and Vivek Raghavan. 2022. Code switched and code mixed speech recognition for indic languages. *arXiv preprint arXiv:2203.16578*.

HoLam Chung, Junan Li, Pengfei Liu, Wai-Kim Leung, Xixin Wu, and Helen Meng. 2022. Improving rare words recognition through homophone extension and unified writing for low-resource cantonese speech recognition. In *2022 13th International Symposium on Chinese Spoken Language Processing (ISCSLP)*, pages 26–30. IEEE.

Anuj Diwan, Rakesh Vaideeswaran, Sanket Shah, Ankita Singh, Srinivasa Raghavan, Shreya Khare, Vinit Unni, Saurabh Vyas, Akash Rajpuria, Chiranjeevi Yarra, et al. 2021. Multilingual and code-switching asr challenges for low resource indian languages. *arXiv preprint arXiv:2104.00235*.

Abhimanyu Dubey, Abhinav Jauhri, Abhinav Pandey, Abhishek Kadian, Ahmad Al-Dahle, Aiesha Letman, Akhil Mathur, Alan Schelten, Amy Yang, Angela Fan, et al. 2024. The llama 3 herd of models. *arXiv preprint arXiv:2407.21783*.

Daya Guo, Dejian Yang, Haowei Zhang, Junxiao Song, Ruoyu Zhang, Runxin Xu, Qihao Zhu, Shirong Ma, Peiyi Wang, Xiao Bi, et al. 2025. Deepseek-r1: Incentivizing reasoning capability in llms via reinforcement learning. *arXiv preprint arXiv:2501.12948*.

Daniil Gurgurov, Tanja Bäumel, and Tatiana Anikina. 2024. Multilingual large language models and curse of multilinguality. *arXiv preprint arXiv:2406.10602*.

Wei-Ning Hsu, Benjamin Bolte, Yao-Hung Hubert Tsai, Kushal Lakhotia, Ruslan Salakhutdinov, and Abdelrahman Mohamed. 2021. Hubert: Self-supervised speech representation learning by masked prediction of hidden units.

Olga Iakovenko and Thomas Hain. 2024. Methods for automatic matrix language determination of code-switched speech. *arXiv preprint arXiv:2410.02521*.

Albert Q Jiang, Alexandre Sablayrolles, Arthur Mensch, Chris Bamford, Devendra Singh Chaplot, Diego de las Casas, Florian Bressand, Gianna Lengyel, Guillaume Lample, Lucile Saulnier, et al. 2023. Mistral 7b. *arXiv preprint arXiv:2310.06825*.

Hexin Liu, Leibny Paola Garcia, Xiangyu Zhang, Andy WH Khong, and Sanjeev Khudanpur. 2024a. Enhancing code-switching speech recognition with interactive language biases. In *ICASSP 2024-2024 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, pages 10886–10890. IEEE.

Hexin Liu, Haihua Xu, Leibny Paola Garcia, Andy WH Khong, Yi He, and Sanjeev Khudanpur. 2023a. Reducing language confusion for code-switching speech recognition with token-level language diarization. In *ICASSP 2023-2023 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, pages 1–5. IEEE.

Hexin Liu, Xiangyu Zhang, Haoyang Zhang, Leibny Paola Garcia, Andy WH Khong, Eng Siong Chng, and Shinji Watanabe. 2024b. Aligning speech to languages to enhance code-switching speech recognition. *arXiv preprint arXiv:2403.05887*.

Pengfei Liu, Weizhe Yuan, Jinlan Fu, Zhengbao Jiang, Hiroaki Hayashi, and Graham Neubig. 2023b. Pre-train, prompt, and predict: A systematic survey of prompting methods in natural language processing. *ACM Computing Surveys*, 55(9):1–35.

Vijayaditya Peddinti, Daniel Povey, and Sanjeev Khudanpur. 2015. A time delay neural network architecture for efficient modeling of long temporal contexts. In *Interspeech*, volume 2015, pages 3214–3218.

Puyuan Peng, Brian Yan, Shinji Watanabe, and David Harwath. 2023. Prompting the hidden talent of web-scale speech models for zero-shot task generalization. *arXiv preprint arXiv:2305.11095*.

Rohit Prabhavalkar, Takaaki Hori, Tara N Sainath, Ralf Schlüter, and Shinji Watanabe. 2023. End-to-end speech recognition: A survey. *IEEE/ACM Transactions on Audio, Speech, and Language Processing*, 32:325–351.

Alec Radford, Jong Wook Kim, Tao Xu, Greg Brockman, Christine McLeavey, and Ilya Sutskever. 2023. Robust speech recognition via large-scale weak supervision. In *International conference on machine learning*, pages 28492–28518. PMLR.

Alec Radford, Jeffrey Wu, Rewon Child, David Luan, Dario Amodei, Ilya Sutskever, et al. 2019. Language models are unsupervised multitask learners. *OpenAI blog*, 1(8):9.

Budi Setiawan. 2023. Code-mixing vs code-switching: a study of grammatical perspective through code-switching varieties. *KnE Social Sciences*, pages 47–57.

Tongtong Song, Qiang Xu, Meng Ge, Longbiao Wang, Hao Shi, Yongjie Lv, Yuqin Lin, and Jianwu Dang. 2022. Language-specific characteristic assistance for code-switching speech recognition. *arXiv preprint arXiv:2206.14580*.

Brij Mohan Lal Srivastava and Sunayana Sitaram. 2018. Homophone identification and merging for code-switched speech recognition. In *Interspeech*, pages 1943–1947.

Jiwon Suh, Injae Na, and Woohwan Jung. 2024. Improving domain-specific asr with llm-generated contextual descriptions. In *Interspeech 2024*, pages 1255–1259.

S Thara and Prabaharan Poornachandran. 2018. Code-mixing: A brief survey. In *2018 International conference on advances in computing, communications and informatics (ICACCI)*, pages 2382–2388. IEEE.

Qinyi Wang, Emre Yılmaz, Adem Derinel, and Haizhou Li. 2019. Code-switching detection using asr-generated language posteriors. *arXiv preprint arXiv:1906.08003*.

Shinji Watanabe, Takaaki Hori, Suyoun Kim, John R Hershey, and Tomoki Hayashi. 2017. Hybrid ctc/attention architecture for end-to-end speech recognition. *IEEE Journal of Selected Topics in Signal Processing*, 11(8):1240–1253.

Genta Indra Winata, Alham Fikri Aji, Zheng-Xin Yong, and Thamar Solorio. 2022. The decades progress on code-switching research in nlp: A systematic survey on trends and challenges. *arXiv preprint arXiv:2212.09660*.

An Yang, Baosong Yang, Beichen Zhang, Binyuan Hui, Bo Zheng, Bowen Yu, Chengyuan Li, Dayiheng Liu, Fei Huang, Haoran Wei, et al. 2024a. Qwen2. 5 technical report. *arXiv preprint arXiv:2412.15115*.

Hao Yang, Jinming Zhao, Gholamreza Haffari, and Ehsan Shareghi. 2023. Investigating pre-trained audio encoders in the low-resource condition. *arXiv preprint arXiv:2305.17733*.

Tzu-Ting Yang, Hsin-Wei Wang, Yi-Cheng Wang, Chi-Han Lin, and Berlin Chen. 2024b. An effective mixture-of-experts approach for code-switching speech recognition leveraging encoder disentanglement. In *ICASSP 2024-2024 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, pages 11226–11230. IEEE.

Yuhang Yang, Yizhou Peng, Hao Huang, Eng Siong Chng, and Xionghu Zhong. 2024c. Adapting openai's whisper for speech recognition on code-switch mandarin-english seame and asru2019 datasets. In *2024 Asia Pacific Signal and Information Processing Association Annual Summit and Conference (APSIPA ASC)*, pages 1–6. IEEE.

Shuaishuai Ye, Shunfei Chen, Xinhui Hu, and Xinkang Xu. 2024. Sc-moe: Switch conformer mixture of experts for unified streaming and non-streaming code-switching asr. *arXiv preprint arXiv:2406.18021*.

Tengfei Yu, Xuebo Liu, Liang Ding, Kehai Chen, Dacheng Tao, and Min Zhang. 2024. Speech sense disambiguation: Tackling homophone ambiguity in end-to-end speech translation. In *Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 8020–8035.

Shuai Zhang, Jiangyan Yi, Zhengkun Tian, Jianhua Tao, Yu Ting Yeung, and Liqun Deng. 2022. Reducing language context confusion for end-to-end code-switching automatic speech recognition. *arXiv preprint arXiv:2201.12155*.

# Beyond Labeled Datasets: Advancing TTS with Direct Preference Optimization on Unlabeled Speech Dataset

**Andrii Zhuravlov[1], Volodymyr Sydorskyi[1],**

[1]National Technical University of Ukraine "Igor Sikorsky Kyiv Polytechnic Institute"

zhuravlov.andrii@lll.kpi.ua, v.sydorskyi@kpi.ua

## Abstract

In our work, we enhance language model-based Text-to-Speech (TTS) training from unlabeled speech data using Direct Preference Optimization (DPO). Given the critical challenges related to the quality and quantity of data required for high-quality speech generation systems, it is essential to develop cost-effective approaches to training such models. We propose a two-stage fine-tuning approach, which extends traditional fine-tuning on texts generated by automatic speech recognition (ASR) models and incorporates direct preference optimization (DPO) along with dataset expansion using texts generated by large language models (LLMs). Experiments and comparisons conducted on two different datasets demonstrate that our approach achieves results comparable to traditional fine-tuning on human-labeled data. The code is publicly available on GitHub[1].

## 1 Introduction

In recent years, the quality of speech generation has significantly improved, largely due to advancements in high-quality audio quantizers such as Hi-FiCodec (Yang et al., 2023) and VQ-VAE, which was used in xTTS system (Casanova et al., 2024). These developments have enabled the use of Transformer architectures (Vaswani, 2017), which are known to perform well with large-scale datasets but are prone to overfitting on smaller datasets.

As a result, data collection and the quality of datasets remain critical challenges in the continued advancement of TTS models. One approach to increasing data availability involves using ASR models to automatically annotate raw audio data. However, this method compromises the quality of speech generation, as raw audio data is often of low quality and ASR models introduce recognition errors. To address these issues, a WV-MOS-based filtering method (Ogun et al., 2023) has been

proposed to improve data set quality by filtering low-quality samples using WV-MOS models. Additionally, it has been demonstrated that raw audio data quality can be improved using noise-filtering systems (Ni et al., 2023; Hao et al., 2021), which boosts TTS model performance but complicates the preprocessing pipeline.

Fortunately, the Transformer architecture enables the application of techniques from NLP, particularly training pipelines for large language models (LLMs). Training LLMs generally consists of two stages: pretraining and fine-tuning. During pretraining, the model is trained on a large corpus of low-quality data, and in the fine-tuning stage, methods such as Reinforcement Learning from Human Feedback (RLHF) (Ouyang et al., 2022) and instruction tuning are used to improve the quality of model outputs. Tian et al. (2024) demonstrated a similar approach and showed that using the Direct Preference Optimization (DPO) algorithm is effective for TTS models. In Tian et al. (2024), preference alignment was guided by three metrics – WER (Whisper-large model (Radford et al., 2023)), SPK_SIM (RawNet (Jung et al., 2024)), and Proxy_MOS (UTMOS (Saeki et al., 2022)) – to evaluate preferences between sample pairs and there was shown a great boost on each metric. Moreover, was shown that improvement does not depend on the exact metric models. In another study (Hussain et al., 2025), preference pairs were constructed using the character error rate (CER) and cosine similarity (SSIM) metrics, along with a modified version of the DPO method – e Reward-aware Preference Optimization (RPO) – to enable more fine-grained preference calibration. The study also demonstrated that employing DPO or RPO for fine-tuning TTS models can lead to improvements in the overall quality of the resulting system.

Inspired by the previous findings (Tian et al., 2024; Hussain et al., 2025), in this work, we

---

[1]https://github.com/BirdWithDreams/beyond-labeled-datasets-tts

present a semi-supervised training strategy for a language model-based text-to-speech (TTS) system, aiming to reduce reliance on labeled data while maintaining high synthesis quality. We explored six different model training strategies on two distinct datasets. Additionally, we proposed a semi-supervised two-stage training strategy: first, a standard fine-tuning on ASR-labeled data, followed by DPO fine-tuning on a combination of original dataset ASR texts and LLM-generated texts.

Our results demonstrate that the proposed training strategy outperforms traditional fine-tuning on human-labeled data in two out of the three primary evaluation metrics. Additionally, we perform a human evaluation using the Comparative Mean Opinion Score (CMOS) methodology. The results indicate that the proposed approach is statistically comparable to conventional fine-tuning techniques on human labeled data. Furthermore, we adapted the popular xTTS framework to support training using the DPO method.

## 2 Method

### 2.1 Semi-supervised training methodology

Since collecting and annotating data for training a TTS model is a complex and resource-intensive process, we developed a fully unsupervised training method. The core idea of our approach consists of two stages: first, a standard fine-tuning of the base model on an ASR-labeled dataset, and second, the creation of a DPO dataset using this model.

Creating the DPO dataset does not require human involvement. All we need is a model checkpoint $M$ for generation, a set of reference audios $A$, and a set of texts $T$. The latter initially consisted of ASR-labeled texts from the original datasets, which we further expanded with LLM-generated texts. The generation procedure is detailed in Appendix A.1. Then, for each $(a, t) \subset A \times T$ pair, we generated 10 audio $\{y_{a,t,1}, y_{a,t,2}, \ldots, y_{a,t,10}\}$ variants using the model $M$. This set of samples was ranked within each pair using three primary evaluation metrics, and the final ranking was determined using a harmonic mean aggregation, sorting the generated audio $y$ from best to worst $(y_{a,t,1}^{f} \succ y_{a,t,2}^{f} \succ \ldots \succ y_{a,t,10}^{f})$ and based on this ranking, win-lose pairs were selected for preference alignment. The full procedure of DPO dataset construction is described in Appendix A.2.



Figure 1: Schematic diagram of our model training pipeline.

### 2.2 Optimization Objectives

In our work, we used two approaches for training models (Figure 1). The first is standard fine-tuning as it was done in xTTS model (Casanova et al., 2024) with cross-entropy loss. Given an input text $x$ (represented as a sequence of tokens), reference audio $r$, and target audio $y$ (represented as a sequence of audio tokens), the model is trained to minimize the following loss function (a full derivation and detailed breakdown of its components is provided in Appendix B):

$$\mathcal{L}_{\text{FN}}(\pi_\theta) = \mathcal{L}_{\text{audio}}(\pi_\theta) + \alpha \cdot \mathcal{L}_{\text{text}}(\pi_\theta) \quad (1)$$

The second approach is Direct Preference Optimization, DPO (Rafailov et al., 2024). A win-lose pair dataset $(r, x, y_w, y_l)$ was constructed, where the sequence $y_w$ is preferred over $y_l$. To optimize the model on such data, we can use the $\mathcal{L}_{\text{DPO}}$ loss as it was described in the original paper (Rafailov et al., 2024).

Instead of using the pure $\mathcal{L}_{\text{DPO}}$ loss, we incorporated an additional $\mathcal{L}_{\text{text}}$ term, as it is known that TTS models perform better when optimized not only with respect to audio but also with respect to the given text itself. Finally, our DPO loss takes the form:

$$\mathcal{L}(\pi_\theta; \pi_{\text{ref}}) = \mathcal{L}_{\text{DPO}}^{\text{audio}}(\pi_\theta; \pi_{\text{ref}}) + \mathcal{L}_{\text{text}}(\pi_\theta) \quad (2)$$

## 3 Experiment Setup

### 3.1 Model

We chose xTTSv2 (Casanova et al., 2024) as our base model. We use it because, firstly, it has an LM-based architecture, which is critically important for DPO fine-tuning. Secondly, it achieves state-of-the-art zero-shot performance in multiple

languages, including English. Additionally, it is highly stable during training, which is crucial for the reproducibility and reliability of our results.

## 3.2 Data

In our work, we used two datasets: The LJ Speech Dataset (Ito and Johnson, 2017) (denoted as $D_{\text{LJ}}^{\text{original}}$), representing a classic single-speaker audiobook-like dataset, and the CSTR VCTK Corpus (Veaux et al., 2017) (denoted as $D_{\text{VCTK}}^{\text{original}}$) as a multi-speaker dataset. Since these datasets contain manually annotated transcriptions, whereas we needed synthetic ones, we generated new transcriptions using Whisper-medium (Radford et al., 2023). We then split the data into a training set and a hold-out set in a ratio 80/20 for The LJ Speech Dataset and 90/10 for CSTR VCTK Corpus, resulting in the following two datasets: $D_{\text{LJ}}^{\text{ASR}}$ and $D_{\text{VCTK}}^{\text{ASR}}$.

Additionally, we generated 15,000 texts using Llama 3.2 3b (Dubey et al., 2024). During generation, we employed a specialized text attribute combinator (considering factors such as length, topic, domain, complexity, etc.) to ensure maximum diversity in the generated texts (see Fig. 2). These texts were later used to augment the original datasets, enhancing their variability and robustness.

## 3.3 Metrics for DPO dataset

To evaluate the models and construct the DPO dataset, we used three main metrics: intelligibility (WER), speaker similarity (SS), and Proxy MOS (PMOS). The following models were used to calculate these metrics: Whisper-Medium (Radford et al., 2023) for WER, ECAPA2 Speaker Embedding Extractor (Thienpondt and Demuynck, 2023) for SS, and UTMOS (Saeki et al., 2022) for PMOS. Model validation was performed on the holdout subsets of our $D_{\text{LJ}}^{\text{original}}$ and $D_{\text{VCTK}}^{\text{original}}$ datasets.

## 3.4 Experiments

For each dataset group, $D_{\text{LJ}}$ and $D_{\text{VCTK}}$, the following fine-tuning experiments were conducted:

1. Fine-tuning (FN) of the base xTTSv2 model on $D^{\text{original}}$.

2. Fine-tuning (FN) of the base xTTSv2 model on $D^{\text{ASR}}$.

3. DPO fine-tuning of the base xTTSv2 model on $D^{\text{DPO}}$.

4. DPO fine-tuning of the base xTTSv2 model on $D^{\text{DPO}} + D^{\text{Generated}}$.

5. DPO fine-tuning of the model from the corresponding checkpoint (LJ-ASR or VCTK-ASR) on $D^{\text{DPO}}$.

6. DPO fine-tuning of the model from the corresponding checkpoint (LJ-ASR or VCTK-ASR) on $D^{\text{DPO}} + D^{\text{Generated}}$.

Let's call these experiments L1-6 for $D_{\text{LJ}}$ group of datasets and V1-6 for $D_{\text{VCTK}}$ group. The validation results for each setup are presented in Table 1 and 2 for the $D_{\text{LJ}}$ and $D_{\text{VCTK}}$ dataset groups, respectively.

Table 1: Model Performance on $D_{\text{LJ}}$ dataset group

| Model | WER ↓ | SS ↑ | PMOS ↑ |
|---|---|---|---|
| Base xTTSv2 | $0.071 \pm 0.008$ | $0.423 \pm 0.003$ | $3.68 \pm 0.016$ |
| L1 (Original) | $0.056 \pm 0.014$ | $\mathbf{0.481 \pm 0.003}$ | $3.816 \pm 0.013$ |
| L2 (ASR) | $0.064 \pm 0.010$ | $0.478 \pm 0.003$ | $3.79 \pm 0.013$ |
| L3 (DPO) | $\mathbf{0.043 \pm 0.003}$ | $0.445 \pm 0.003$ | $3.733 \pm 0.012$ |
| L4 (DPO) | $0.064 \pm 0.011$ | $0.465 \pm 0.002$ | $\mathbf{3.959 \pm 0.010}$ |
| L5 (DPO) | $0.110 \pm 0.012$ | $0.432 \pm 0.003$ | $2.821 \pm 0.011$ |
| L6 (DPO) | $0.224 \pm 0.035$ | $0.417 \pm 0.003$ | $2.392 \pm 0.012$ |

Table 2: Model Performance on $D_{\text{VCTK}}$ dataset group

| Model | WER ↓ | SS ↑ | PMOS ↑ |
|---|---|---|---|
| Base xTTSv2 | $0.020 \pm 0.004$ | $0.481 \pm 0.014$ | $3.895 \pm 0.026$ |
| V1 (Original) | $0.041 \pm 0.007$ | $\mathbf{0.500 \pm 0.014}$ | $3.685 \pm 0.029$ |
| V2 (ASR) | $0.055 \pm 0.009$ | $0.494 \pm 0.014$ | $3.630 \pm 0.030$ |
| V3 (DPO) | $0.014 \pm 0.003$ | $0.471 \pm 0.015$ | $4.009 \pm 0.022$ |
| V4 (DPO) | $\mathbf{0.013 \pm 0.003}$ | $0.482 \pm 0.016$ | $\mathbf{4.108 \pm 0.019}$ |
| V5 (DPO) | $0.273 \pm 0.037$ | $0.412 \pm 0.014$ | $2.662 \pm 0.047$ |
| V6 (DPO) | $0.087 \pm 0.013$ | $0.453 \pm 0.014$ | $3.324 \pm 0.043$ |

## 3.5 Fine-Tuning vs. DPO Fine-Tuning

When comparing these two training methods, the first noticeable trend is that standard fine-tuning achieves the best SS metric across both datasets: 0.481 for $D_{\text{LJ}}$ (L1) and 0.5 for $D_{\text{VCTK}}$ (V1). As expected, training on ASR-labeled data (L2, V2) performs worse than training on human-labeled data (L1, V1). However, DPO training on ASR-labeled data (L3, V3) either outperforms or at least matches traditional fine-tuning with a cross-entropy objective on human-labeled data (L1, V1).

Interestingly, the best results in speech naturalness (PMOS metric) are achieved when the dataset is expanded with LLM-generated data (L4, V4), even surpassing a PMOS score of 4. Regarding intelligibility (WER metric), the best performance in the $D_{\text{VCTK}}$ dataset group (WER 0.013) is also obtained with DPO tuning on the expanded

dataset (`V4`), outperforming both classical fine-tuning on human-labeled data (`V1`, WER 0.04) and the xTTSv2 baseline (WER 0.02).

For the $D_{\text{LJ}}$ dataset group, the best WER score is achieved by the `L3` model (WER 0.041), outperforming both `L1` (WER 0.053) and the baseline (WER 0.071). However, the DPO fine-tune on the expanded dataset (`L4`) achieves results similar to traditional fine-tuning on human-labeled data and worse than DPO tune on unexpanded dataset (`L3`). This behavior can be explained by the nature of the The LJ Speech Dataset (Ito and Johnson, 2017). This dataset consists of audiobook recordings where audio is sometimes segmented inaccurately, resulting in partial sentences, such as only the beginning or end of a sentence like "`According to Secretary Dillon,`" or "`iron and the like in combination with phosphoric, sulphuric and other acids.`". Expanding the $D_{\text{LJ}}$ dataset with LLM-generated texts, which consist of fully formed sentences, does not necessarily improve model performance within the $D_{\text{LJ}}$ dataset.

In contrast, the CSTR VCTK Corpus (Veaux et al., 2017) dataset, which was created by reading newspaper sentences rather than slicing pre-existing audio, is more aligned with the way LLM-generated texts are structured. This explains why in the $D_{\text{VCTK}}$ dataset group, fine-tuning on the expanded dataset (`V4`) yields better results (WER 0.013, SS 0.481, PMOS 4.108) than fine-tuning on standard texts (`V3`) (WER 0.014, SS 0.471, PMOS 4.009).

### 3.6 Effects of ASR Checkpoint Initialization

Comparing experiments 5-6 with 3-4, we observe that fine-tuning from the ASR checkpoint consistently yields worse results than fine-tuning from the base model on the same data. `L5` and `V5` show much higher WER (0.109 and 0.269, respectively) and lower PMOS (2.821 and 2.665). This can be explained by the fact that standard fine-tuning narrows the generation space, whereas DPO fine-tuning only adjusts the probability distribution within that space without altering it. In other words, the "softer" DPO fine-tuning from a checkpoint with greater generation variability leads to better results than fine-tuning from a checkpoint with lower variability. This holds true even though, in the latter case, the model was explicitly trained to reproduce the distribution of a specific dataset.

However, we observe that in both cases (`L4`,

`L6` and `V4`, `V6`) adding AI-generated texts improves model performance across almost all metrics (except for WER in the `L3-4` cases), supporting our hypothesis that expanding the dataset with AI-generated texts positively impacts model quality.

### 3.7 CMOS validation

To further evaluate the proposed approach, we conducted a CMOS (Comparative Mean Opinion Score) validation following the methodology detailed in Appendix C.1. CMOS evaluation was conducted on four experimental pairs: (1 vs. 4), (2 vs. 4), (1 vs. 6), and (2 vs. 6). The evaluation considers two criteria: speaker similarity (SS), and a combined metric reflecting both naturalness and intelligibility (CM). Results are presented in Table 3.

Table 3: Method Pair Comparison Data

| Comparison | SS | CM |
|---|---|---|
| Exp. 1 vs Exp. 4 | $-0.12$ | $-0.06$ |
| Exp. 1 vs Exp. 6 | $0.88$ | $0.21$ |
| Exp. 2 vs Exp. 4 | $-0.21$ | $-0.17$ |
| Exp. 2 vs Exp. 6 | $0.75$ | $0.14$ |

Positive values in Table 3 indicate that the first experiment in the pair is preferred over the second, while negative values indicate the opposite.

In the comparison between Exp. 1 and Exp. 4, the SS metric marginally favors Exp. 4 ($-0.12$), while the combined metric indicates near equivalence ($-0.06$). Similarly, for the Exp. 2 vs. Exp. 4 comparison, both metrics slightly favor Exp. 4, with $-0.21$ for SS and $-0.17$ for CM.

More substantial differences are observed with Exp. 6. In both (1 vs. 6) and (2 vs. 6) comparisons, the metrics are positive (e.g., 0.88 and 0.75 for SS and combined in 1 vs. 6), indicating that the classical tuning baselines were preferred. These results support our conclusions based on automatic validation metrics (WER, SS, PMOS). Additionally, we performed a statistical significance analysis of the results, detailed in Appendix C.2. We also provide a detailed subgroup CMOS analysis in Appendix C.3. Overall, the CMOS evaluation indicates that the proposed method, particularly in Exp. 4, achieves quality comparable to conventional fine-tuning using human-labeled data.

## 4 Conclusion

We have developed a two-stage training strategy for TTS models based on DPO fine-tuning. We proposed a fully unsupervised training pipeline for TTS models and demonstrated that it can achieve results comparable to traditional supervised fine-tuning on human-labeled data. This approach significantly reduces costs, as manual annotation requires substantial resources and time. Therefore, our method is more efficient without sacrificing model quality.

Additionally, we showed that expanding original datasets with LLM-generated texts substantially improves the naturalness (PMOS) of generated audio while having a mixed impact on intelligibility (WER), which requires further investigation across different data types and datasets.

## 5 Limitations

Given the limitations of our work, we used the high-quality xTTSv2 model as our baseline. For future research, it would be valuable to train several models from scratch – one on ASR-labeled data and one on human-labeled data and compare how DPO fine-tuning affects their quality. Another interesting direction is to compare our method of constructing win-lose pairs for DPO with human-based pair selection.

Our pipeline involves components that may introduce or amplify societal biases:

1. ASR-Induced Bias: we rely on an Automatic Speech Recognition (ASR) model (Whisper-medium) to generate transcripts for unlabeled audio. It is well-documented that ASR systems can have higher error rates for speakers with non-native accents, certain dialects, or speech impediments. Such transcription errors may degrade the quality of synthesized speech for already underrepresented groups, potentially reinforcing existing biases in the system.

2. LLM-Induced Bias: The use of a Large Language Model (Llama 3) to generate supplementary text for training introduces the risk of inheriting its intrinsic biases. While we employed an attribute combinator to encourage text diversity (Appendix A.1), the generated content may still reflect dominant cultural viewpoints or stereotypes present in the LLM's training data.

Future work should involve auditing the model's performance across more diverse demographic groups and developing methods to mitigate any identified biases.

## 6 Ethical concerns

The development of advanced Text-to-Speech (TTS) technologies, such as the one presented in this paper, carries significant societal implications that warrant careful consideration. We are committed to the responsible advancement of AI and outline the primary ethical concerns related to our work below.

The most significant risk associated with high-fidelity TTS is the potential for misuse in creating synthetic audio, often referred to as "deepfakes."

- Using unauthorized voice synthesis to impersonate someone for fraudulent purposes, such as deceiving individuals or bypassing voice authentication systems.

- Disinformation and propaganda involve fabricating audio evidence to spread misinformation, defame individuals, or manipulate public opinion.

- Generating non-consensual audio content to harass or bully.

While our research aims to advance machine learning methodology, we recognize this dual-use nature. We advocate for the development and adoption of robust safeguards, such as audio watermarking techniques and detection models for synthetic speech, which should accompany any deployment of this technology in real-world applications.

this tool, the authors reviewed and edited the content as needed and take full responsibility for the publication's content.

# References

Edresson Casanova, Kelly Davis, Eren Gölge, Görkem Göknar, Iulian Gulea, Logan Hart, Aya Aljafari, Joshua Meyer, Reuben Morais, Samuel Olayemi, et al. 2024. Xtts: a massively multilingual zero-shot text-to-speech model. *arXiv preprint arXiv:2406.04904*.

Abhimanyu Dubey, Abhinav Jauhri, Abhinav Pandey, Abhishek Kadian, Ahmad Al-Dahle, Aiesha Letman, Akhil Mathur, Alan Schelten, Amy Yang, Angela Fan, et al. 2024. The llama 3 herd of models. *arXiv preprint arXiv:2407.21783*.

Xiang Hao, Xiangdong Su, Radu Horaud, and Xiaofei Li. 2021. Fullsubnet: A full-band and sub-band fusion model for real-time single-channel speech enhancement. In *ICASSP 2021-2021 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, pages 6633–6637. IEEE.

Shehzeen Hussain, Paarth Neekhara, Xuesong Yang, Edresson Casanova, Subhankar Ghosh, Mikyas T Desta, Roy Fejgin, Rafael Valle, and Jason Li. 2025. Koel-tts: Enhancing llm based speech generation with preference alignment and classifier free guidance. *arXiv preprint arXiv:2502.05236*.

Keith Ito and Linda Johnson. 2017. The lj speech dataset. https://keithito.com/LJ-Speech-Dataset/.

Jee-weon Jung, Wangyou Zhang, Jiatong Shi, Zakaria Aldeneh, Takuya Higuchi, Barry-John Theobald, Ahmed Hussen Abdelaziz, and Shinji Watanabe. 2024. Espnet-spk: full pipeline speaker embedding toolkit with reproducible recipes, self-supervised front-ends, and off-the-shelf models. *arXiv preprint arXiv:2401.17230*.

Zhaoheng Ni, Sravya Popuri, Ning Dong, Kohei Saijo, Xiaohui Zhang, Gael Le Lan, Yangyang Shi, Vikas Chandra, and Changhan Wang. 2023. Exploring speech enhancement for low-resource speech synthesis. *arXiv preprint arXiv:2309.10795*.

Sewade Ogun, Vincent Colotte, and Emmanuel Vincent. 2023. Can we use common voice to train a multi-speaker tts system? In *2022 IEEE Spoken Language Technology Workshop (SLT)*, pages 900–905. IEEE.

Long Ouyang, Jeffrey Wu, Xu Jiang, Diogo Almeida, Carroll Wainwright, Pamela Mishkin, Chong Zhang, Sandhini Agarwal, Katarina Slama, Alex Ray, et al. 2022. Training language models to follow instructions with human feedback. *Advances in neural information processing systems*, 35:27730–27744.

Alec Radford, Jong Wook Kim, Tao Xu, Greg Brockman, Christine McLeavey, and Ilya Sutskever. 2023. Robust speech recognition via large-scale weak supervision. In *International conference on machine learning*, pages 28492–28518. PMLR.

Rafael Rafailov, Archit Sharma, Eric Mitchell, Christopher D Manning, Stefano Ermon, and Chelsea Finn. 2024. Direct preference optimization: Your language model is secretly a reward model. *Advances in Neural Information Processing Systems*, 36.

Takaaki Saeki, Detai Xin, Wataru Nakata, Tomoki Koriyama, Shinnosuke Takamichi, and Hiroshi Saruwatari. 2022. Utmos: Utokyo-sarulab system for voicemos challenge 2022. *arXiv preprint arXiv:2204.02152*.

Jenthe Thienpondt and Kris Demuynck. 2023. Ecapa2: A hybrid neural network architecture and training strategy for robust speaker embeddings. In *2023 IEEE Automatic Speech Recognition and Understanding Workshop (ASRU)*, pages 1–8. IEEE.

Jinchuan Tian, Chunlei Zhang, Jiatong Shi, Hao Zhang, Jianwei Yu, Shinji Watanabe, and Dong Yu. 2024. Preference alignment improves language model-based tts. *arXiv preprint arXiv:2409.12403*.

A Vaswani. 2017. Attention is all you need. *Advances in Neural Information Processing Systems*.

Christophe Veaux, Junichi Yamagishi, Kirsten MacDonald, et al. 2017. Cstr vctk corpus: English multi-speaker corpus for cstr voice cloning toolkit. *University of Edinburgh. The Centre for Speech Technology Research (CSTR)*, 6:15.

Dongchao Yang, Songxiang Liu, Rongjie Huang, Jinchuan Tian, Chao Weng, and Yuexian Zou. 2023. Hifi-codec: Group-residual vector quantization for high fidelity audio codec. *arXiv preprint arXiv:2305.02765*.

## A  Experiment Details

### A.1  Generation of LLM texts

To enhance the diversity of LLM-generated texts, we used a specialized attribute combinator to construct prompts for the LLM. Fig. 2 shows what attributes, sub-attributes and constraints were used to create high variability of generated texts.

All of this enables the creation of a highly diverse vocabulary, addressing one of the key challenges in TTS model training – bias between over-represented and underrepresented words. At the same time, using LLM-generated texts helps to fill in gaps in the vocabulary and allows for fine control over text types and formats. This can be particularly useful when fine-tuning a model for a highly specialized domain with limited original data.

### A.2  Creating DPO datasets

In the first stage, the LJ-ASR and VCTK-ASR models were trained using standard fine-tuning on ASR-generated texts from the base xTTSv2 model. Then, using the latest checkpoints of these models, the datasets $D_{\text{LJ}}^{\text{DPO}}$, $D_{\text{LJ}}^{\text{Generated}}$, $D_{\text{VCTK}}^{\text{DPO}}$, and $D_{\text{VCTK}}^{\text{Generated}}$ were constructed.

**Method for Constructing $D_{\text{LJ}}^{\text{Generated}}$:** a selection of audio samples was taken from the original dataset $D_{\text{LJ}}^{original}$, and AI-generated texts were evenly distributed among them. For each (audio, text) pair, 10 samples $(y_{a,t,1}, y_{a,t,2}, \dots, y_{a,t,10})$ were generated using the LJ-ASR model, and evaluation metrics were computed for each sample using our evaluation models. Notice, that each $y$ is not a generated audio, but a sequence of audio codes produced by LM head (see Casanova et al. (2024)).

Next, these samples were ranked from best to worst according to each metric $(y_{a,t,1}^{\text{wer}} \succ y_{a,t,2}^{\text{wer}} \succ \dots \succ y_{a,t,10}^{\text{wer}})$. Based on their ranking, a normalized score between 0 and 1 was assigned to each sample.

$$metric\_rank = place/10$$

To determine the final ranking, we calculated harmonic mean of our metrics' ranks:

$$f\_rank = \frac{3}{\frac{1}{wer\_rank} + \frac{1}{ss\_rank} + \frac{1}{mos\_rank}}.$$

Then, based on its values, the preferred ($y_w$) and less preferred ($y_l$) samples were selected. We choose them as the second sample from each

edge, mean $y_w = y_{a,t,2}^{\text{f\_rank}}$ and $y_l = y_{a,t,9}^{\text{f\_rank}}$. The most extreme samples, the absolute best ($y_1^{\text{f\_rank}}$) and worst ($y_{10}^{\text{f\_rank}}$), were excluded to ensure that the preference optimization for the model was not overly obvious. Following this process, the $D_{\text{LJ}}^{\text{Generated}}$ dataset was constructed: $(a, t, y_w, y_l)$, where $a$ is reference audio sample, $t$ - reference text, $y_w$ - preferred sequence of audio codes and $y_l$ - non-preferred sequence of audio codes.

**Method for Constructing $D_{\text{VCTK}}^{\text{Generated}}$:** Since $D_{\text{VCTK}}^{\text{original}}$ contains 108 unique speakers and 13,000 unique texts—where different speakers may read the same text—the dataset includes a total of 44,000 (speaker, text) pairs. Each speaker has between 200 and 500 recordings. We decided to construct the DPO version of this dataset in a similar manner. Our 15,000 LLM-generated texts were evenly distributed among all speakers, with repetitions, ensuring that each speaker had an average of 500 unique texts. The subsequent sample generation, ranking, and win-lose pair selection followed the same approach as for $D_{\text{LJ}}^{\text{Generated}}$, with the VCTK-ASR model used during sample generation.

**Construction of $D_{\text{LJ}}^{\text{DPO}}$ and $D_{\text{VCTK}}^{\text{DPO}}$:** The datasets $D_{\text{LJ}}^{\text{DPO}}$ and $D_{\text{VCTK}}^{\text{DPO}}$ were constructed similarly to $D_{\text{LJ}}^{\text{Generated}}$ and $D_{\text{VCTK}}^{\text{Generated}}$, with the key difference that the texts were taken from the original $D_{\text{LJ}}^{\text{original}}$ and $D_{\text{VCTK}}^{\text{original}}$ datasets.

## B  Objectives definitions

Classical cross-entropy (CE) loss on text and audio tokens:

$$\begin{aligned}
\mathcal{L}_{\text{FN}}(\pi_\theta) &= \mathcal{L}_{\text{audio}}(\pi_\theta) + \alpha \cdot \mathcal{L}_{\text{text}}(\pi_\theta) \\
&= -\mathbb{E}_{(x,r,y)\sim\mathcal{D}} \log \pi_\theta(y \mid x, r) \\
&\quad - \alpha \cdot \mathbb{E}_{(x,r,y)\sim\mathcal{D}} \log \pi_\theta(x^t \mid x^{t-1}, r)
\end{aligned} \tag{3}$$

DPO loss from the original paper (Rafailov et al., 2024):

$$\mathcal{L}_{\text{DPO}}(\pi_\theta; \pi_{\text{ref}}) =$$
$$= -\mathbb{E}\left[ \log \sigma \left( \beta \log \frac{\pi_\theta(y_w \mid x)}{\pi_{\text{ref}}(y_w \mid x)} - \beta \log \frac{\pi_\theta(y_l \mid x)}{\pi_{\text{ref}}(y_l \mid x)} \right) \right] \tag{4}$$

where $\pi_\theta$ is the model that is being optimized and $\pi_{\text{ref}}$ is the original model.

Final objective for second stage of proposed method:

Figure 2: Combinator's attributes diagram

$$\mathcal{L}\left(\pi_\theta ; \pi_{\mathrm{ref}}\right)=\mathcal{L}_{\mathrm{DPO}}^{\mathrm{audio}}\left(\pi_\theta ; \pi_{\mathrm{ref}}\right)+\mathcal{L}_{\mathrm{text}}\left(\pi_\theta\right)$$

$$=-\mathbb{E}\left[\log \sigma\left(\beta \log \frac{\pi_\theta\left(y_w \mid x\right)}{\pi_{\mathrm{ref}}\left(y_w \mid x\right)}-\beta \log \frac{\pi_\theta\left(y_l \mid x\right)}{\pi_{\mathrm{ref}}\left(y_l \mid x\right)}\right)\right]$$

$$-\alpha \mathbb{E} \log \pi_\theta\left(x^t \mid x^{t-1}, r\right) \tag{5}$$

## C  CMOS validation details

### C.1  Validation methodology

To facilitate the CMOS validation process, an automated service was developed to efficiently and conveniently collect user feedback. This service operates in a fully automated mode and presents evaluation tasks in a user-friendly format. The selection of speakers and texts for CMOS validation was intentionally made diverse: from the VCTK dataset, two speakers (one male and one female) were chosen for each of five distinct accents—American, British, Indian, Irish, and Scottish—resulting in a total of ten speakers, with an additional speaker selected from the LJ-Speech dataset. For each speaker, four different texts were selected from the validation subsets of the original datasets ($D_{\mathrm{LJ}}^{\mathrm{original}}$ and $D_{\mathrm{VCTK}}^{\mathrm{original}}$), and synthetic audio was generated for these texts using the corresponding TTS model.

The survey methodology is as follows: for each participant, the service randomly selects three

speakers from the ten available in $D_{\mathrm{VCTK}}^{\mathrm{original}}$ and adds one speaker from $D_{\mathrm{LJ}}^{\mathrm{original}}$. The participant is then presented with 32 evaluation items: four method pairs × four speakers × four comparisons per speaker. Each evaluation item is structured as follows: "Please assess which of the two audio samples better corresponds to the reference recording according to a specific criterion, using a scale from $-3$ to 3, where 3 indicates that the first sample is significantly better, 0 means both are approximately equal, and $-3$ indicates that the second sample is significantly better."

### C.2  Statistical significance analysis

To assess the reliability of the CMOS evaluation, we conducted a statistical significance test by evaluating the null hypothesis that the mean CMOS score is zero. A $p$-value above the significance threshold (0.05) indicates that the compared models are statistically equivalent, whereas a value below the threshold suggests a significant preference for one model over the other. The outcomes of this analysis for the SS and CM metrics are presented in Table 4 and Table 5, respectively.

The statistical analysis of the CMOS metrics (Table 4 and Table 5) reveals that Experiment 6, which implements the proposed method, significantly underperforms in speaker similarity compared to tra-

356

Table 4: Statistical significance of the SS metric

|  | 1–4 | 1–6 | 2–4 | 2–6 |
|---|---|---|---|---|
| Mean | −0.12 | 0.88 | −0.21 | 0.75 |
| $t$-test $p$-val | 0.4 | **0.001** | 0.22 | **0.007** |
| Wilcoxon $p$-val | 0.5 | **0.001** | 0.25 | **0.015** |
| Sample size | 83 | 80 | 76 | 73 |

Table 5: Statistical significance of the CM metric

|  | 1–4 | 1–6 | 2–4 | 2–6 |
|---|---|---|---|---|
| Mean | −0.06 | 0.21 | −0.17 | 0.14 |
| $t$-test $p$-val | 0.74 | 0.34 | 0.39 | 0.52 |
| Wilcoxon $p$-val | 0.74 | 0.34 | 0.35 | 0.49 |
| Sample size | 83 | 80 | 76 | 73 |

ditional approaches. In particular, Experiments 1 and 2 show strong and statistically significant advantages over Experiment 6 in the SS metric, with $p$-values well below 0.001. In contrast, comparisons involving Experiment 4 do not exhibit significant differences, indicating that this configuration achieves perceptual speaker similarity comparable to traditional fine-tuning.

For the Combined Metric (CM), none of the comparisons across experimental conditions yield statistically significant differences (all $p > 0.3$), suggesting that all methods perform similarly in terms of overall speech naturalness and text accuracy. These findings indicate that the proposed method, especially in Experiment 4, maintains competitive perceptual quality, while Experiment 6 demonstrates limited effectiveness in preserving vocal identity.

## C.3 CMOS on specific groups

To further investigate the behavior of the proposed models across different speaker characteristics, we conducted a stratified CMOS analysis by accent and gender. As described in Appendix C.1, we examined the same four experimental method pairs: (1 vs. 4), (2 vs. 4), (1 vs. 6), and (2 vs. 6). For each pair, two criteria were evaluated: speaker similarity (SS) and a composite metric capturing clarity, naturalness, and intelligibility (CM). Results are presented in Table 6 and Table 7.

Overall, positive values in Tables 6 and 7 indicate a preference for the first method in each comparison, while negative values indicate preference for the second.

The proposed Method 4 (DPO training with generated data) demonstrates advantages in perceptual quality (CM) across several accents, with moderate gains in speaker similarity (SS). Compared

Table 6: CMOS Results by Accent

| Method Pair | Accent | SS | CM | N |
|---|---|---|---|---|
| 1–4 | American | −0.545 | 0.818 | 11 |
|  | English | −0.235 | −0.706 | 17 |
|  | Indian | 0.300 | 0.500 | 10 |
|  | Irish | −0.533 | −0.400 | 15 |
|  | Scottish | 0.444 | −0.111 | 9 |
|  | lj | 0.048 | −0.095 | 21 |
| 1–6 | American | 1.300 | 0.800 | 5 |
|  | English | 1.400 | −1.200 | 5 |
|  | Indian | 0.933 | 0.333 | 15 |
|  | Irish | 1.231 | 0.000 | 13 |
|  | Scottish | 0.400 | −0.050 | 20 |
|  | lj | 1.000 | 0.647 | 17 |
| 2–4 | American | 0.833 | −0.333 | 6 |
|  | English | −0.600 | 0.000 | 10 |
|  | Indian | −0.263 | −0.526 | 19 |
|  | Irish | −0.333 | 0.000 | 12 |
|  | Scottish | 0.545 | 0.364 | 11 |
|  | lj | −0.667 | −0.278 | 18 |
| 2–6 | American | 1.417 | 0.417 | 12 |
|  | English | 0.909 | 0.091 | 11 |
|  | Indian | −0.071 | −0.214 | 14 |
|  | Irish | 0.667 | 0.111 | 9 |
|  | Scottish | 0.444 | 0.222 | 9 |
|  | lj | 1.056 | 0.222 | 18 |

to the baseline trained on human-annotated data (Method 1), Method 4 achieves better CM scores for the LJ speaker (−0.095), English (−0.706), and Irish (−0.400), and also shows improved SS for English (−0.235) and Irish (−0.533), suggesting enhanced or preserved speaker identity. Relative to the ASR-supervised baseline (Method 2), Method 4 again receives more favorable CM values for American-accented speech (−0.333) and Scottish (−0.364), along with strong SS improvements for English (−0.600) and the LJ speaker (−0.667), highlighting its robustness on several accent groups. However, performance on some accents, such as Indian and American in the 1–4 comparison, remains challenging. By contrast, Method 6 (DPO with generated data initialized from a pretrained checkpoint) mostly underperforms relative to both baseline methods across individual accent groups, showing less consistent gains in either CM or SS. It is also important to note the variability in group sizes ($N$), with some accent groups containing relatively few samples. This limits the statistical robustness of per-accent conclusions and calls for caution when interpreting fine-grained differences.

Gender-based analysis further supports the effectiveness of the proposed Method 4, particularly for female speakers. Compared to the human-

Table 7: CMOS Results by Gender

| Method Pair | Gender | SS | CM | N |
|---|---|---|---|---|
| 1–4 | F | −0.292 | −0.250 | 24 |
| | M | −0.105 | 0.079 | 38 |
| | lj (F) | 0.048 | −0.095 | 21 |
| 1–6 | F | 0.656 | 0.031 | 32 |
| | M | 1.032 | 0.161 | 31 |
| | lj (F) | 1.000 | 0.647 | 17 |
| 2–4 | F | −0.379 | −0.069 | 29 |
| | M | 0.241 | −0.207 | 29 |
| | lj (F) | −0.667 | −0.278 | 18 |
| 2–6 | F | 0.300 | 0.150 | 30 |
| | M | 0.857 | 0.086 | 35 |
| | lj (F) | 1.056 | 0.222 | 18 |

annotated baseline (Method 1), Method 4 achieves better CM scores for female speakers (−0.250) and the LJ speaker (−0.095), while also improving SS for females (−0.292), indicating that the proposed approach is preferred in terms of both perceptual quality and speaker similarity. For male speakers, results are more mixed: while SS is slightly better (−0.105), the CM score (0.079) indicates a mild preference for the baseline. In comparison to the ASR-supervised baseline (Method 2), Method 4 again shows lower CM for female (−0.069) and male (−0.207) speakers, and achieves a strong improvement for the LJ speaker (−0.278), with consistent SS gains for females (−0.379) and the LJ speaker (−0.667), reinforcing the robustness of Method 4 for female voices. In contrast, Method 6 performs worse than both baselines across all gender groups. CM scores are consistently positive when compared to both Method 1 and Method 2, indicating that listeners preferred the baseline systems in terms of clarity, naturalness, and intelligibility. SS values also show degradation, with all comparisons yielding positive scores, suggesting less accurate speaker identity preservation. As with the accent-based analysis, these observations should be interpreted with caution due to relatively small group sizes ($N$), especially for the LJ speaker.

To complement the CMOS evaluation, we conducted statistical significance testing on the speaker similarity (SS) and clarity/naturalness (CM) scores within each subgroup. For every experimental method pair and demographic subgroup (by gender and accent), we applied one-sample t-tests and Wilcoxon signed-rank tests against a null hypothesis of zero (i.e., no perceived difference between systems). The resulting p-values are presented in

Table 8 (gender) and Table 9 (accent).

The statistical significance analysis supports the earlier observations (Table 4 and Table 5). For Method 4, p-values across most gender groups are above the 0.05 threshold in both SS and CM comparisons against baseline Methods 1 and 2, indicating no statistically significant difference and suggesting that the proposed method performs comparably to the baselines. In contrast, Method 6 consistently shows statistically significant differences in SS when compared to both baselines (e.g., $p < 0.05$ for both males and females), pointing to a degradation in speaker similarity. For CM, however, most comparisons yield p-values above 0.05, implying that the perceptual quality of Method 6 is not significantly different from the baselines despite the SS drop.

For accent-based comparisons, the majority of p-values also exceed 0.05, which may reflect a lack of statistical power due to small sample sizes within individual accent groups. Nevertheless, a few accents (e.g., Irish and American in 1–6 and 2–6 pairs) show marginal or significant effects, particularly in SS, indicating that accent-specific behavior may warrant closer examination in future studies with larger cohorts.

These findings underscore the importance of subgroup-level analysis in evaluating TTS systems. Listener demographics—such as gender and accent—can influence judgments of speaker similarity and perceptual quality, and adequate subgroup representation is crucial to draw robust, generalizable conclusions.

Table 8: Statistical Significance of CMOS Scores by Gender

| Method Pair | Gender | SS $t$-p | SS $w$-p | CM $t$-p | CM $w$-p |
|---|---|---|---|---|---|
| 1–4 | M | 0.612 | 0.618 | 0.761 | 0.766 |
| | F | 0.307 | 0.349 | 0.434 | 0.532 |
| | lj (F) | 0.867 | 0.769 | 0.820 | 0.744 |
| 1–6 | M | 0.0001 | 0.0006 | 0.643 | 0.603 |
| | F | 0.0001 | 0.0006 | 0.662 | 0.606 |
| | lj (F) | 0.063 | 0.078 | 0.287 | 0.223 |
| 2–4 | M | 0.452 | 0.504 | 0.546 | 0.552 |
| | F | 0.163 | 0.189 | 0.828 | 0.729 |
| | lj (F) | 0.014 | 0.023 | 0.508 | 0.520 |
| 2–6 | M | 0.0001 | 0.0006 | 0.782 | 0.776 |
| | F | 0.410 | 0.392 | 0.679 | 0.622 |
| | lj (F) | 0.015 | 0.023 | 0.664 | 0.668 |

Table 9: Statistical Significance of CMOS Scores by Accent

| Method Pair | Accent | SS $t$-p | SS $w$-p | CM $t$-p | CM $w$-p |
|---|---|---|---|---|---|
| 1–4 | Scottish | 0.498 | 0.531 | 0.834 | 1.000 |
| | Indian | 0.343 | 0.531 | 0.363 | 0.424 |
| | Irish | 0.072 | 0.072 | 0.361 | 0.404 |
| | lj | 0.867 | 0.769 | 0.820 | 0.744 |
| | American | 0.216 | 0.030 | 0.156 | 0.055 |
| | English | 0.431 | 0.463 | 0.055 | 0.054 |
| 1–6 | Scottish | 0.237 | 0.227 | 0.878 | 0.975 |
| | Indian | 0.025 | 0.034 | 0.559 | 0.523 |
| | Irish | 0.009 | 0.004 | 1.000 | 1.000 |
| | lj | 0.063 | 0.078 | 0.287 | 0.223 |
| | American | 0.004 | 0.008 | 0.236 | 0.254 |
| | English | 0.478 | 0.750 | 0.109 | 0.188 |
| 2–4 | Scottish | 0.327 | 0.336 | 0.420 | 0.539 |
| | Indian | 0.426 | 0.365 | 0.213 | 0.169 |
| | Irish | 0.529 | 0.624 | 1.000 | 1.000 |
| | lj | 0.014 | 0.023 | 0.508 | 0.520 |
| | American | 0.259 | 0.375 | 0.679 | 0.750 |
| | English | 0.193 | 0.219 | 1.000 | 1.000 |
| 2–6 | Scottish | 0.447 | 0.516 | 0.708 | 0.844 |
| | Indian | 0.856 | 0.917 | 0.609 | 0.667 |
| | Irish | 0.169 | 0.250 | 0.824 | 1.000 |
| | lj | 0.015 | 0.023 | 0.664 | 0.668 |
| | American | 0.002 | 0.008 | 0.499 | 0.550 |
| | English | 0.074 | 0.110 | 0.884 | 0.902 |



Figure 3: DPO training loss

# D  Additional Experimental Results

## D.1  DPO optimzation

Figure 3 illustrates the average training loss across several model variants using DPO. Most models demonstrate a smooth and consistent decrease in loss, indicating stable convergence behavior. While some variance exists across configurations, there are no signs of divergence or abrupt fluctuations. Overall, these results suggest that training with DPO is stable under the tested conditions.

# L1RA: Dynamic Rank Assignment in LoRA Fine-Tuning

**Raul Singh**[*], **Nicolò Brunello**[*], **Vincenzo Scotti**[†] and **Mark James Carman**[*]

[*]DEIB, Politecnico di Milano
Via Ponzio 34/5, 20133, Milano (MI), Italy
[†]KASTEL, Karlsruhe Institute of Technology (KIT)
Am Fasanengarten 5, 76131, Karlsruhe, Germany
raul.singh@mail.polimi.it     nicolo.brunello@polimi.it
vincenzo.scotti@kit.edu     mark.carman@polimi.it

## Abstract

The ability of Large Language Models (LLMs) to solve complex tasks has made them crucial in the development of AI-based applications. However, the high computational requirements to fine-tune these LLMs on downstream tasks pose significant challenges, particularly when resources are limited. In response to this challenge, we introduce L1RA, a novel technique aimed at dynamically distributing the rank of low-rank adapters during fine-tuning using LoRA. Given a rank budget (i.e., total sum of adapters rank), L1RA leverages $L_1$ regularisation to prune redundant ranks and redistribute them across adapters, thereby optimising resource utilisation. Through a series of comprehensive experiments, we empirically demonstrate that L1RA maintains comparable or even reduced computational overhead compared to other LoRA variants, including the vanilla approach, while achieving same or better performances. Moreover, the post-training analysis of rank distribution unveiled insights into the specific model components requiring the most adaptation to align with the task objective: the feed-forward layers and the attention output projection. These results highlight the efficacy of L1RA in not only enhancing the efficiency of LLM fine-tuning, but also in providing valuable diagnostic information for model refinement and customisation. In conclusion, L1RA stands as a promising technique for advancing the performance and interpretability of LLM adaptation, particularly in scenarios where computational resources are constrained.

## 1 Introduction

*Large Language Models* (LLMs) have revolutionised *Natural Language Processing* (NLP) and *Artificial Intelligence* (AI) (Zhao et al., 2023), enabling sophisticated applications. LLM's language understanding and generation capabilities make them suitable for an impressive number of applications (Raffel et al., 2020; Brown et al., 2020;

Sanh et al., 2022). Moreover, their adoption as core for *chatbots* (Scotti et al., 2024) have made them essential for the final consumers of this technology. However, to excel in these specific tasks, even conversation, LLMs often require *fine-tuning*, a process essential for tailoring their vast pre-trained knowledge to new specific contexts and domains. This adaptation ensures optimal performance and task alignment, making fine-tuning a critical step in deploying LLMs effectively.

The fine-tuning process, however, presents challenges, particularly concerning computational resources. Adaptation to specific domains, such as chatbot dialogue or instruction-following tasks, demands substantial computational power, which may be impractical or unfeasible in resource-constrained environments. Recent advancements in efficient fine-tuning techniques, including *Low-Rank Adaptation* (LoRA) (Hu et al., 2022), *prefix tuning* (Li and Liang, 2021) and the *gradient-free methods* like *Memory-efficient Zeroth-order Optimiser* (MeZO) (Malladi et al., 2023), offer promising solutions. These techniques leverage strategies like low-rank parameterisation to reduce computational overhead, making fine-tuning more accessible.

In this paper, we introduce $L_1$-*regularised Rank Assignment* (L1RA): a technique aimed at enhancing the efficiency and effectiveness of LLM fine-tuning. L1RA extends LoRA by introducing $L_1$ regularisation to enforce rank sparsity and dynamic rank allocation during training to get the best from the available resources. Assuming a given *rank budget* (i.e., total sum of LoRA adapter ranks), L1RA prunes redundant ranks and reallocates them across adapters during the fine-tuning process. We pair L1RA with our tool *Memory GPU Estimation of LLM Allocation for Training Optimisation* (MEMORY-GELATO) to be sure to match available resources constraints. Through a series of experiments, ranging from small-scale analyses to comprehensive comparisons

with other fine-tuning techniques, we evaluate the performance of L1RA. The results highlight how L1RA can offer better comparable results to alternative LORA variants reallocating ranks with negligible difference in resources consumption and better results even with respect to regular LORA.

We divide the rest of the paper into the following sections. In Section 2, we present the related works on efficient LLM fine-tuning. In Section 3, we explain the reasons behind our work. In Sections 4 and 5, we describe, respectively, the L1RA fine-tuning algorithm and the MEMORY-GELATO tool. In Section 6, we outline the experiments to evaluate our model and in Section 7 we present the obtained results. In Section 8 we comment on the results we obtained. Finally, in Section 9, we sum up our work and suggest possible future extensions.

## 2   Related works

Efficient fine-tuning techniques have garnered increasing attention lately, due to the computational demands associated with adapting LLMs to specific tasks. The proposed techniques evolved significantly during the last few years. Initial approaches like *Transformer Adapters* (Houlsby et al., 2019; Bapna and Firat, 2019) introduced additional parameters in the form of a pair of linear projections with a bottleneck in the middle, increasing network depth and latency, thereby hindering scalability. In response, LORA-based solutions (Hu et al., 2022) have emerged as a promising alternative. LORA addresses the limitations of adapters by introducing low-rank parameterisation, effectively reducing the number of parameters needed for adaptation. This technique has gained widespread adoption for its ability to achieve efficient fine-tuning without compromising performance. Alternative techniques like MEZO (Malladi et al., 2023) target the training algorithm rather than the network structure, focusing on fine-tuning through forward passes only, eliminating the need for backpropagation and the subsequent overhead. Other approaches like prefix-tuning (Li and Liang, 2021) learn only the embeddings of a *continuos prompt* that can be used as a prefix to the input to condition the LLM output towards the desired task. Among these techniques, LORA stands out as the most adopted due to its effectiveness in balancing computational efficiency, performance and ease of use.

As premised, LORA operates by introducing pairs of low-rank matrices $\mathbf{A} \in \mathbb{R}^{d_{in} \times r}$ and $\mathbf{B} \in$



Figure 1: LORA adapters: pre-trained weights are frozen while the two adapter matrices are updated during the fine-tuning.

$\mathbb{R}^{r \times d_{out}}$ into the network architecture (see Figure 1); the product $\Delta \mathbf{W} \in \mathbb{R}^{d_{in} \times d_{out}}$ of these two matrices encodes the weights difference induced by fine-tuning for a specific weight matrix $\mathbf{W} \in \mathbb{R}^{d_{in} \times d_{out}}$ of the pre-trained model, as explained by Equation (1). During fine-tuning, these adapter matrices are updated while the original LLM parameters are kept frozen. By leveraging low-rank parameterisation, LORA effectively reduces the computational overhead associated with fine-tuning while preserving the expressive power of the LLM. Moreover, this approach has demonstrated empirically significant improvements in efficiency without sacrificing performance across various downstream tasks.

$$\mathbf{h} = \mathbf{x} \cdot (\mathbf{W} + \Delta \mathbf{W}) = \mathbf{x} \cdot \mathbf{W} + \mathbf{x} \cdot (\mathbf{A} \cdot \mathbf{B}) \quad (1)$$

While LORA offers notable benefits, several variants have been proposed to address specific limitations or further enhance its capabilities. Examples of these alternative solutions are those aimed at stabilising the training process, like LORA+ (Hayou et al., 2024), which introduces a matrix-specific scaling parameter on the learning rate to improve performances and convergence time, and *Rank-Stabilised* LORA (RSLORA) (Kalajdzievski, 2023), which uses a rank correcting factor to prevent gradient collapse. Other variants, like *Quantised* LORA (QLORA) (Dettmers et al., 2023), aim at further reducing computational complexity by heavily quantising the base model weights (for reduced memory requirements and increased inference speed) while operating on floating-point representation of the trainable weights (for numeric precision and, thus, training stability), thereby improving the overall efficiency. In this paper we focus on techniques for

adaptive rank allocation. In fact, LoRA adapters depend on the rank hyper-parameter, which can be selected dynamically for each pair of adapter matrices. In this sense, some solutions have been proposed to tackle the issue of rank selection in order to (i) get rid of unused parameters and (ii) find the best possible rank allocation allowed by the available memory.

One of the first solutions for dynamic rank allocation was presented with ADALoRA (Zhang et al., 2023), which enforces a *Singular Value Decomposition*-inspired (SVD-inspired) decomposition of the adapter weights through additional regularisation terms in the loss. Further refinements of this technique came with *Sparse* LoRA (SoRA) (Ding et al., 2023), which uses an intermediate gating mechanism with $L_1$ regularisation and *proximal gradient descent* to iteratively reduce the used ranks, and, *Vector-based Random matrix Adaptation* (VeRA) (Kopiczko et al., 2023), which reduces the trainable LoRA parameters through shared random weights matrices and works on rank allocation updating only layer-specific parameter vectors. In parallel, *Dynamic rank selection* LoRA (DyLoRA) (Valipour et al., 2023), proposed a solution exploring a range of possible ranks during training to find the optimal ones for each matrix.

## 3 Motivations

LoRA adapters represent a valuable step towards end-user fine-tuning of LLMs, making this technology more accessible and customisable. The existence of techniques like ADALoRA, SoRA or DyLoRA allowing for dynamic rank and pruning (i.e., removing the $i$-th column in $\mathbf{A}$ and the $i$-th row in $\mathbf{B}$) are the results of advances towards better exploitation of computational resources. Hereafter, we highlight some points of improvement for ADALoRA and SoRA (the main solutions for dynamic rank allocation), in terms of computational resources exploitation, that are motivating our work.

ADALoRA proposes a SVD-inspired formulation of the adapter:

$$\Delta \mathbf{W} = \mathbf{U} \cdot \boldsymbol{\Sigma} \cdot \mathbf{V}^\top = \mathbf{U} \cdot \mathrm{diag}(\boldsymbol{\sigma}) \cdot \mathbf{V}^\top \quad (2)$$

where $\mathbf{U} \in \mathbb{R}^{(d_{in} \times r)}$, $\mathbf{V} \in \mathbb{R}^{(d_{out} \times r)}$, $\boldsymbol{\sigma} \in \mathbb{R}_0^{+r}$. Then, it enforces an additional regularisation term $\mathcal{L}_{SVD}(\Delta \mathbf{W})$ to the loss to imposing orthonormality on the adapter matrices.

$$\Delta \mathbf{W} = \|\mathbf{U}^\top \cdot \mathbf{U} - \mathbf{I}\|_2^2 + \|\mathbf{V}^\top \cdot \mathbf{V} - \mathbf{I}\|_2^2 \quad (3)$$

Despite this constraint allows to interpret the values of $\boldsymbol{\sigma}$ as the eigenvalues and, thus, prune all elements corresponding to null eigenvalues in increases the memory and time requirements of the training process with respect to a normal LoRA.

SoRA builds on top of ADALoRA, discarding the SVD constraint and substituting the vector of eigenvalues with a gating vector $\mathbf{g} \in \mathbb{R}^r$ and enforcing sparsity adding to the loss a L1 regularisation penalty on $g$. This simple, yet effective solution, encourages to prune all elements corresponding to a $0$ valued element in the gate, as they will be ignored in the computation of the output (exactly as the elements corresponding to a null eigenvalue). The complete formulation of SoRA includes the proximal gradient update using a thresholding function that ensures training stability. This addition is already part of the optimiser we use in our experiments (see Appendix B for further details).



Figure 2: Motivating example: $r'$ and $r''$ are such that $r' + r'' = 2r$, so that the total amount of adapters memory is the same with and without optimal allocation.

All the proposed solutions for dynamic rank assignment correctly work to reduce the rank used in the adapter matrices. However, these LoRA variants are limited in the sense that they do not allow for spare (unused) ranks re-assignment and they rather wait for the end of training to prune the matrices. They instead propose starting directly from higher ranks, usually $3r/2$, which increase the overall mem-

ory requirement with respect to a base LORA operating with the same resources and rank $r$. Consider the toy example in Figure 2, where we have the comparison between the matrices of LORA with fixed rank allocation and the matrices with performance-optimal rank allocation. In this case we would have a rank budget of $2r$ that, in the performance-optimal allocation, is divided between $r'$, in the first adapter, and $r''$, in the second adapter, that $r' + r'' = 2r$ and $r' > r > r''$. In this configuration, with adapters like ADALORA or SORA, we would need to start at least from a rank budget of $2r' > 2r$ to reach the performance optimal allocation, which is above the available budget of $2r$. Moreover, it may be the case where, since we are talking of constrained resources, the model with all the adapters starting from rank $r'$ would not fit in memory.

Besides the theoretical aspects of staying within the rank budget, we also have a "physical" constraint given by the amount of available GPU memory. To tackle this problem we developed the MEMORY-GELATO tool, which comes as a complement to L1RA. Though accurate estimates of the memory usage we can identify the starting rank without exceeding the available resources. Similarly to other solutions, L1RA can drop the ranks in excess, but differently from the other takes care of re-allocating at runtime those ranks, all of this staying within the given budget.

In this section we described exactly the problems we tackle with our work: *how to get the best performances given a fixed rank or memory budget*? In other words, our contribution is an algorithm that dynamically re-allocates rank amongst adapter matrices in order to maximise performance given a fixed maximum memory budget available, complemented with a tool for memory budget estimation.

## 4 L1RA



$$\mathbf{h} \in \mathbb{R}^d$$

$$\oplus$$

$$\mathbf{W} \in \mathbb{R}^{d \times d}$$

$$\mathbf{B} \in \mathbb{R}^{r \times d}$$
$$b_{i,j} = 0$$

$$\mathbf{c} \in \mathbb{R}^r$$
$$c_i = 1$$

$$\mathbf{A} \in \mathbb{R}^{d \times r}$$
$$a_{i,j} \sim \mathcal{N}(0, \sigma^2)$$

Pre-trained weights    L1RA Adapter

$$\mathbf{x} \in \mathbb{R}^d$$

Figure 3: L1RA adapters

L1RA adapters, depicted in Figure 3, extend the LORA framework by introducing rank pruning and reallocation mechanisms within a fixed rank or memory budget. The goal of L1RA is to identify the performance-optimal rank configuration in computational constrained settings where memory –and time– may be limited. This dynamic rank adjustment ensures that the model efficiently utilises the available resources, enhancing performance without exceeding the same constraints a vanilla LORA adapter would have.

Mathematically, given an input vector $\mathbf{x} \in \mathbb{R}^{d_{in}}$, we compute the output $\mathbf{h} \in \mathbb{R}^{d_{out}}$ of a L1RA adapter as described in Equation (4). Where $\mathbf{W} \in \mathbb{R}^{d_{in} \times d_{out}}$ is the original matrix of pre-trained weights, $\Delta \mathbf{W} \in \mathbb{R}^{d_{in} \times d_{out}}$ is the adapter matrix of weights decomposed in $\mathbf{A} \in \mathbb{R}^{d_{in} \times r}$, $\mathbf{c} \in \mathbb{R}^r$ and $\mathbf{B} \in \mathbb{R}^{r \times d_{out}}$, and the $\mathrm{diag}(\cdot)$ function outputs a diagonal matrix with the elements of the input vector as values on the diagonal. The rank $r$ depends on the specific adapter and is selected through optimisation during training.

$$\mathbf{h} = \mathbf{x} \cdot (\mathbf{W} + \Delta \mathbf{W}) = \mathbf{x} \cdot \mathbf{W} + \mathbf{x} \cdot (\mathbf{A} \cdot \mathrm{diag}(\mathbf{c}) \cdot \mathbf{B}) \tag{4}$$

The $\mathbf{c}$ vector we introduced, similar to the gating system of SORA, is a technical device to ease the enforcing of the sparsity constraint. We could have obtained the same effect of imposing sparsity on $\mathbf{c}$ (resulting from the regularisation term $\lambda \|\mathbf{c}\|_1$), by applying the same $L_1$ constraint to the columns of the input projection matrix $\mathbf{A}$ of a regular LORA adapter. Doing so would have resulted, however, in much slower redistribution of rank, since an entire column of $\mathbf{A}$ would need to have converged to zero before it could be removed and reassigned to a different matrix, whereas a single component of the $\mathbf{c}$ vector falling to zero is sufficient for reassignment. Thus, while the $\mathbf{c}$ vector does introduce a small number of additional parameters, it results in faster and more direct rank-sparsification, while not affecting the overall transformation of the adapter.

We report the training process of a model using L1RA adapters in the pseudocode detailed in Algorithm 1 and we detail the rank pruning and re-allocation process in Figure 4. The overall training is similar to that of a model using LORA adapters. The loss function is changed to include the $L_1$ regularisation term (controlled by the $\lambda$ hyperparameter) on the elements of the $\mathbf{c}$ vector. Similarly to SORA, by enforcing sparsity on the $\mathbf{c}$ vector through this regularisation, we achieve rank pruning. In fact, whenever an element of $\mathbf{c}$ is shrunk

**Algorithm 1** L1RA pseudocode

**Require:**
- $\vartheta$ ▷ Model parameters
- $\mathcal{D}$ ▷ Data
- $r \in \mathbf{N}^+$ ▷ Initial adapters rank

$\Delta\vartheta \leftarrow \{\}$ ▷ Adapter parameters
**for** $\mathbf{W} \in \vartheta$ **do** ▷ Initialise adapters of all layers
  $\mathbf{A} \leftarrow \mathbf{A} \in \mathbb{R}^{d \times r} \sim \mathcal{N}(0, \sigma^2)$
  $\mathbf{B} \leftarrow \mathbf{0} \in \{0\}^{r \times d}$
  $\mathbf{c} \leftarrow \mathbf{1} \in \{1\}^r$
  $\Delta\vartheta \leftarrow \Delta\vartheta \cup \{(\mathbf{A}, \mathbf{B}, \mathbf{c})\}$
**end for**
**for** $i \in [0, n_{epochs}) \subseteq \mathbb{N}$ **do** ▷ Iterate over epochs
  **for** $X \in \mathcal{D}$ **do** ▷ Iterate over training samples
    $\mathcal{L}(\Delta\vartheta) \leftarrow - \ln\ P(X; \vartheta, \Delta\vartheta) + \lambda \cdot \sum_{(\mathbf{A}, \mathbf{B}, \mathbf{c}) \in \Delta\vartheta} \|\mathbf{c}\|_1$ ▷ Get loss
    $\Delta\vartheta \leftarrow \Delta\vartheta - \eta \cdot \nabla_{\Delta\vartheta} \cdot \mathcal{L}(\Delta\vartheta)$ ▷ Update adapter weights
    $\rho \leftarrow 0$ ▷ Initialise spare ranks
    $\Delta\vartheta_u \leftarrow []$ ▷ Initialise list of unpruned adapters
    **for** $(\mathbf{A}, \mathbf{B}, \mathbf{c}) \in \Delta\vartheta$ **do** ▷ Iterate over adapters
      **if** $\exists c \in \mathbf{c} | c = 0$ **then** ▷ Check for a rank decrease
        $\rho \leftarrow \rho + \sum_{c \in \mathbf{c}} \mathcal{I}(c = 0)$ ▷ Count spare ranks
        $(\mathbf{A}, \mathbf{B}, \mathbf{c}) \leftarrow f_{prune}(\mathbf{A}, \mathbf{B}, \mathbf{c})$ ▷ Apply pruning
      **else** ▷ Else if not pruned
        $f_{insert}(\Delta\vartheta_u, (\mathbf{A}, \mathbf{B}, \mathbf{c}))$ ▷ Save adapter for reallocation
      **end if**
    **end for**
    **while** $\rho > 0$ **do** ▷ While there are spare ranks
      **for** $(\mathbf{A}, \mathbf{B}, \mathbf{c}) \in \Delta\vartheta_u$ **do** ▷ Iterate over unpruned adapters
        **if** $\rho > 0$ **then** ▷ if there are spare ranks
          $(\mathbf{A}, \mathbf{B}, \mathbf{c}) \leftarrow f_{reallocate}(\mathbf{A}, \mathbf{B}, \mathbf{c})$ ▷ Re-allocate a rank
          $\mathbf{c} \leftarrow \mathbf{c} / \sum_{c \in \mathbf{c}} c$ ▷ Normalise **c** vector
          $\rho \leftarrow \rho - 1$ ▷ Update spare ranks
        **end if**
      **end for**
    **end while**
  **end for**
**end for**
**return** $\Delta\vartheta$



(a) When a null component in the **c** vector of an adapter is detected, the corresponding elements of the adapter are removed using the $f_{prune}(\cdot)$ function, generating a spare rank that will be reallocate.



(b) When a spare rank is available and a needs to be reallocated, the elements on the target adapter are extended by the $f_{reallocate}(\cdot)$ function (values are initialised as in a regular initialisation).

Figure 4: L1RA pruning and reallocation (lighter colours are for lower absolute values, white is 0).

to 0 the corresponding column in $\mathbf{A}$ and row in $\mathbf{B}$ –the other matrices of the adapter– can be dropped (this is the role of the $f_{prune}(\cdot)$ function).

All the spare ranks generated by this pruning process can be re-allocated to the other, unpruned, adapters. Whenever spare ranks are available, the algorithm cycles over the unpruned adapters *sorted by decreasing order of the minimum value in the* **c** *vector*, so that

$$(\mathbf{A}_i, \mathbf{B}_i, \mathbf{c}_i) > (\mathbf{A}_j, \mathbf{B}_j, \mathbf{c}_j) \iff \min \mathbf{c}_i > \min \mathbf{c}_j \tag{5}$$

and re-assigns a rank to each adapter until spare ranks are no longer available. In other words, each available additional rank is always redistributed to the particular adapter which is most in need of the rank increase, because its current rank-budget

is in full use, with the various components of the **c** vector furthest from zero.

The ordering of unpruned adapters is performed in Algorithm 1 by the $f_{insert}(\cdot)$ function when saving them in $\Delta\vartheta_u$. After the rank re-allocation step, the training procedure reprises. This sorting step is inspired by SVD: assuming that in high-dimensional space the matrices $\mathbf{A}$ and $\mathbf{B}$ can be treated as orthogonal and the **c** vector mimics the diagonal of the singular value matrix.

Compared to other dynamic rank adapters like ADALORA, SORA and DYLORA, L1RA offers significant advantages. If we consider a model using vanilla LORA adapters with a given rank $r$, since all these other techniques do not account for spare ranks re-allocation, they would require starting from a higher rank initialisation to have the adapters requiring a rank $r' > r$ reach that value, implicitly requiring more memory than the original LORA would have used. In contrast, L1RA basically maintains almost the same memory usage by reallocating ranks within the fixed budget (as we detail better in

Section 8, it cannot always be the same due to some weight matrices having $d_{in}$ or $d_{out}$ different from others). Additionally, ADALORA increases the requirements on time and memory by imposing SVD behavior to the elements of the adapter through additional terms in the loss function. L1RA's approach avoids these additional constraints, ensuring computational efficiency while achieving performance-optimal rank configuration and maintaining memory limits. This makes L1RA a better choice for resource-constrained environments, offering a balanced solution for dynamic rank adaptation.

## 5  MEMORY-GELATO

The MEMORY-GELATO tool is crucial to reach full memory exploitation. In fact, it provides an accurate estimate of the memory required to train a model We identified the following contribution to memory estimates:

- *Model parameters*, which include the weights of all layers and the adapters and is influenced by the numeric precision and quantisation;

- *Steady state memory*, that is all memory reserved to keep track of the intermediate states generated by passing data through the model, the gradients and the optimiser state;

- *Activation*, that is the additional memory used to memorise the activations for gradient checkpointing (reducing the memory footprint of gradients);

- *Loss*, which includes the output logits and the memory used to compute the negative log-likelihood;

- *Other contributions*, which includes all the additional elements increasing memory, like operations at the end of the forward pass and the beginning of the backward pass.

To assess the goodness of MEMORY-GELATO estimates, we compared the predicted and real values of memory peak usage for different models, different maximum sequence lengths and different batch sizes. In Figure 5, we can see the difference between the estimates and the real values; while, in Table 1, we report quantitative metrics on estimates goodness. Overall, the error in estimated peak memory usage differs from the real one of a few hundreds MBs, including the overestimate we introduced for safety.

Table 1: MEMORY-GELATO performance in predicting peak memory usage (MAE: Mean Absolute Error, $\rho$: Spearman correlation coefficient, $r$ Pearson correlation coefficient).

| Model | MAE [MB] | $\rho$ | $r$ |
|---|---|---|---|
| MISTRAL 7B V0.3 | 203.05 | 1.0000 | 0.9998 |
| LLAMA2 7B | 109.80 | 1.0000 | 0.9999 |
| LLAMA 3.1 8B | 159.01 | 1.0000 | 0.9999 |
| PHI-3 MINI 4K | 146.03 | 1.0000 | 0.9998 |

## 6  Evaluation

To evaluate L1RA against other adapter approaches, we applied it to fine-tune a LLM in a realistic use case: assistant fine-tuning. Moreover, to demonstrate empirically the practical advantages of L1RA against other approaches we used MEMORY-GELATO to configure the experiment to maximise memory utilisation. We detail the experimental settings in Appendix B.

We experimented fine-tuning to different LLMs (namely MISTRAL 7B V0.3 (Jiang et al., 2023) and LLAMA 3.1 8B (Dubey et al., 2024), both quantised at 4 bits precision) to make sure that L1RA is agnostic of the LLM. We selected the OPENORCA data set (Mukherjee et al., 2023)[1] for this assistant fine-tuning.

In this experiment, we compared the test performance and resource consumption of L1RA against LORA, ADALORA. We compared against two versions of ADALORA: one targeting the same average rank as LORA and L1RA starting from an higher rank (1.5 times that of LORA as suggested in the ADALORA documentation), and another version starting from the same rank of LORA and L1RA and targeting a smaller rank (so that the initial one was 1.5 times that of LORA, again, as suggested in the ADALORA documentation). This fine-tuning task was chosen to demonstrate the practical application of L1RA in efficient fine-tuning of LLMs, particularly in scenarios where fine-tuning on consumer-level GPUs is challenging (e.g., when we reach the limit of usable memory).

Throughout the experiment, we kept track of ranks evolution to analyse the final distribution at the end of training. It this way we can get a better understanding of which components within the Transformer architecture need a more precise

---

[1]Data set card: `https://huggingface.co/datasets/Open-Orca/OpenOrca`

Figure 5: Comparison of peak memory usage estimates from MEMORY-GELATO against actual peak memory usage during training with LORA adapters.

adaptation (identified as those with a higher adapter rank) and shed lights on the internal mechanisms of the Transformer architecture.

# 7 Results



(a) LLAMA 3.1 8B.



(b) MISTRAL 7B V0.3.

Figure 6: Matrix-wise evolution of layer-wise average L1RA adapters rank during training.



Figure 7: Matrix-wise distribution of layer-wise average L1RA adapters rank at the end of training (error bars show standard deviation).



Figure 8: Layer-wise evolution of matrix-wise average L1RA adapters rank at the end of training.

We report the main results of the experiments in Table 2, while the relative values, to ease the comparison, are in Table 3. L1RA achieves the lowest absolute perplexity (PPL) score, improving over both LORA and ADALORA. Moreover, L1RA achieves also the closest training time to that of LORA, with less than 1% difference from LORA. Memory consumption, on the other hand, seems to be similar among the three approaches (most differences from LORA are below 2%) with ADALORA performing better than L1RA (and even LORA in one

configuration). The number of adapter (trainable) parameters shows how ADALORA not applying an actual pruning the matrices and requiring an higher starting rank to target the same average ranks of LORA and L1RA increases significantly the number of parameters (50%) without an improvement on the PPL. On the other side, L1RA exchanging freely parameters between matrices of different sizes causes an increase in the number of parameters as training continues; however, the increase is smaller than that of ADALORA and achieves better PPL than both LORA and ADALORA. We discuss

| Model | Approach | Rank | PPL ↓ | Training time [s] ↓ | Memory [GB][1] ↓ | No. of adapter parameters [M] ↓ | |
|---|---|---|---|---|---|---|---|
| | | | | | | **Start of training** | **End of training** |
| LLAMA 3 8B | LORA | 16 | 3.32 | 30994.89 | 13.84 | 41.94 | 41.94 |
| | ADALORA | 24→16 | 3.63[2] | 32980.28 | 14.00 | 62.92 | 62.92 |
| | ADALORA | 16→12 | 3.57[2] | 32964.14 | 13.76 | 41.95 | 41.95 |
| | L1RA | 16 | **3.25** | 31246.40 | 14.23 | 41.95 | 45.16 |
| MISTRAL 7B V0.3 | LORA | 16 | 2.93 | 37891.88 | 13.58 | 41.94 | 41.94 |
| | ADALORA | 24→16 | 3.16[2] | 40234.87 | 13.82 | 62.92 | 62.92 |
| | ADALORA | 16→12 | 3.16[2] | 40215.02 | 13.59 | 41.95 | 41.95 |
| | L1RA | 16 | **2.91** | 37968.91 | 13.94 | 41.95 | 50.06 |

[1] Values measured using PYTORCH utility for measuring GPU device memory usage: https://pytorch.org/docs/stable/generated/torch.cuda.max_memory_allocated.html.

[2] Values are slightly altered because PPL was computed from the loss of the model which included also the regularisation term, separate computations showed that ADALORA PPL was higher than that of LORA and L1RA.

Table 3: Relative results and resources consumption from Table 2 normalised to the LORA fine-tuning.

| Model | Approach | Rank | Δ PPL [%] ↓ | Δ Training time [%] ↓ | Δ Memory [%] ↓ | Δ No. of adapter parameters [%][1] ↓ |
|---|---|---|---|---|---|---|
| LLAMA 3 8B | ADALORA | 24→16 | 9.34 | 6.41 | 1.16 | 50.02 |
| | ADALORA | 16→12 | 7.53 | 6.35 | −0.58 | 0.02 |
| | L1RA | 16 | −2.11 | 0.81 | 2.82 | 7.68 |
| MISTRAL 7B V0.3 | ADALORA | 24→16 | 7.85 | 6.18 | 1.77 | 50.02 |
| | ADALORA | 16→12 | 7.85 | 6.13 | 0.07 | 0.02 |
| | L1RA | 16 | −0.68 | 0.20 | 2.65 | 19.36 |

[1] Values computed on end-of-training parameters.



(a) LLAMA 3.1 8B start of training.

(b) MISTRAL 7B V0.3 start of training.

(c) LLAMA 3.1 8B halfway through training.

(d) MISTRAL 7B V0.3 halfway through training.

(e) LLAMA 3.1 8B end of training.

(f) MISTRAL 7B V0.3 end of training.

Figure 9: Layer-wise and Matrix-wise evolution of L1RA adapters rank during training.

better about memory consumption and number of adapter parameters in Section 8.

In Figure 6 we can the average evolution of the ranks organised per matrix of the Transformer architecture. As we can see, LLAMA and MISTRAL have the same trends: matrices coming from the

*Feed-Forward Neural Network* (FFNN) layer of the Transformer architecture (up-projection $\mathbf{W}_{up}$, gate $\mathbf{W}_{gate}$, and down-projection $\mathbf{W}_{down}$) are more "rank hungry" than those of the *multi-head self-attention* layers (key $\mathbf{W}_k$, value $\mathbf{W}_v$, query $\mathbf{W}_q$, and output $\mathbf{W}_o$). At the end of training, the difference is clear across all layers, as shown by the averaged rank counts in Figure 7. From Figure 8 we can see another common trend between the two models: ranks are higher in the layers closer to the output of the neural network.

Finally, to report on the individual ranks of each matrix in the Transformer stack, we can see in Figure 9 ranks distributions at the beginning, halfway through and at the end of training. The darker area emerging at the bottom corresponds to the the matrices of the FFNN. We can see how the "rank mass" is higher in these layers especially toward the top of the Transformer network (bottom-right side on the plot) and how it is lower for the multi-head self-attention layer matrices at the bottom of the Transformer network. Moreover, we can wee how with LLAMA this trend is emerging slower: the matrix showing rank distribution at the end is closer to that of MISTRAL halfway through training. Given the higher PPL, we can assume that LLAMA could have been trained for more iterations.

## 8 Discussion

Values of memory consumption does not comply with our expectations, especially if compared taking into account the rank distributions and the number of trainable parameters. Considering the average ranks and the total number of parameters, we expected to see ADALORA starting from the higher rank having the highest memory consumption and ADALORA starting from the same rank as LORA still consume more memory due to the additional operations to compute the regularisation term, while the memory is even lower in the case of LLAMA. We suspect this is due to some internal optimisation or offloading of the trainer in the HUGGINGFACE's TRANSFORMERS library we used (Wolf et al., 2020). Despite we were not able to locate the source of this difference, we conducted a small experiments on the same data using the same handmade training loop with all adapters and we measured an overall higher memory consumption that was more in line with the number of parameters and the compared techniques. In the next iteration of L1RA we plan to drop the trainer to have more reliable estimates.

To comment on the difference in number of parameters between L1RA and LORA, we can see that despite L1RA not exceeding the rank budget, the amount of parameters (and used memory) is slightly higher than LORA. This is a result of allocating the spare ranks to other adapters working on matrices of different sizes. In particular, as we saw from Figure 7, many ranks are allocated to the feed-forward layers, which have a higher ($4\times$) inner projection dimensionality. Despite this situation, L1RA still achieves a lower resources utilisation when compared to ADALORA.

Finally, to comment on the trends observed in Figures 6 to 9, we can say that trends hint how the FFNN layers at the top of the Transformer stack are contributing more to the task being solved. The high-level features processed in that part of the Transformer network need more precise refinement thus the higher rank. Similarly, we believe that exploiting different information from other tokens in the context is not as important as extracting more refined patterns with the non-linear transformations of the FFNN to have the LLM behave as a chatbot assistant, thus the higher ranks in FFNN layers. This observation agrees with intuition that the higher layers of the network should contribute the most to adapting the network to a specific domain, and that the output and FFNN layers are crucial for storing domain-specific information (as noted by (Geva et al., 2021; Biderman et al., 2023)) that likely needs to be updated by the adapters.

## 9 Conclusion

In this paper, we introduced L1RA, a novel technique for efficient LLM fine-tuning. By effectively exploiting the dynamic rank assignment given by $L_1$ regularisation and re-assigning the spare ranks within the available budget, L1RA represents a significant advancement in efficient fine-tuning, offering a promising solution for resource-constrained environments. We completed L1RA with MEMORY-GELATO our tool for GPU memory estimation we can exploit to determine the memory –and thus rank– budget. At this moment we foresee two possible, complementary, directions in the further development of L1RA: we are interested in studying the rank distribution across different models and at a different scales or number of parameter and data-set sizes and we are interested in better understanding the convergence of the proposed method.

## Limitations

In this paper, we mainly focused on the development of L1RA for efficient fine-tuning and its evaluation on realistic use cases, rather than exhaustive experiments. The first limitation is in the choice of the LLM: as for now, we evaluated the results using only MISTRAL 7B v0.3 and LLAMA 3 8B. A proper evaluation would require exploring other openly accessible models of the same and different sizes that would fit on a consumer-level GPU. The second limitation is the choice of the evaluation data set: we considered only the task of instruction fine-tuning since it is a common use case and since it covers many tasks an LLM is required to solve, however a more extensive evaluation exploring different tasks would improve the understanding of L1RA's capabilities.

## Ethics Statement

The authors do not foresee any considerable risks associated with the work presented in this paper. In principle, the L1RA algorithm is intended to make fine-tuning of LLMs more efficient and the MEMORY-GELATO tool is intended for estimating memory consumption such fine-tunings. The authors made the source code publicly available to ensure the reproducibility of the experiments. Refer to Appendix A for further details.

## Acknowledgements

## References

Ankur Bapna and Orhan Firat. 2019. Simple, scalable adaptation for neural machine translation. In *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing, EMNLP-IJCNLP 2019, Hong Kong, China, November 3-7, 2019*, pages 1538–1548. Association for Computational Linguistics.

Stella Biderman, USVSN Sai Prashanth, Lintang Sutawika, Hailey Schoelkopf, Quentin Anthony, Shivanshu Purohit, and Edward Raff. 2023. Emergent and predictable memorization in large language models. In *Advances in Neural Information Processing Systems 36: Annual Conference on Neural Information Processing Systems 2023, NeurIPS 2023, New Orleans, LA, USA, December 10 - 16, 2023*.

Tom B. Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, Sandhini Agarwal, Ariel Herbert-Voss, Gretchen Krueger, Tom Henighan, Rewon Child, Aditya Ramesh, Daniel M. Ziegler, Jeffrey Wu, Clemens Winter, Christopher Hesse, Mark Chen, Eric Sigler, Mateusz Litwin, Scott Gray, Benjamin Chess, Jack Clark, Christopher Berner, Sam McCandlish, Alec Radford, Ilya Sutskever, and Dario Amodei. 2020. Language models are few-shot learners. In *Advances in Neural Information Processing Systems 33: Annual Conference on Neural Information Processing Systems 2020, NeurIPS 2020, December 6-12, 2020, virtual*.

Tri Dao. 2024. Flashattention-2: Faster attention with better parallelism and work partitioning. In *The Twelfth International Conference on Learning Representations, ICLR 2024, Vienna, Austria, May 7-11, 2024*. OpenReview.net.

Tim Dettmers, Artidoro Pagnoni, Ari Holtzman, and Luke Zettlemoyer. 2023. Qlora: Efficient finetuning of quantized llms. In *Advances in Neural Information Processing Systems 36: Annual Conference on Neural Information Processing Systems 2023, NeurIPS 2023, New Orleans, LA, USA, December 10 - 16, 2023*.

Ning Ding, Xingtai Lv, Qiaosen Wang, Yulin Chen, Bowen Zhou, Zhiyuan Liu, and Maosong Sun. 2023. Sparse low-rank adaptation of pre-trained language models. In *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing, EMNLP 2023, Singapore, December 6-10, 2023*, pages 4133–4145. Association for Computational Linguistics.

Abhimanyu Dubey, Abhinav Jauhri, Abhinav Pandey, Abhishek Kadian, Ahmad Al-Dahle, Aiesha Letman, Akhil Mathur, Alan Schelten, Amy Yang, Angela Fan, Anirudh Goyal, Anthony Hartshorn, Aobo Yang, Archi Mitra, Archie Sravankumar, Artem Korenev, Arthur Hinsvark, Arun Rao, Aston Zhang, Aurelien Rodriguez, Austen Gregerson, Ava Spataru, Baptiste Roziere, Bethany Biron, Binh Tang, Bobbie Chern, Charlotte Caucheteux, Chaya Nayak, Chloe Bi, Chris Marra, Chris McConnell, Christian Keller, Christophe Touret, Chunyang Wu, Corinne Wong, Cristian Canton Ferrer, Cyrus Nikolaidis, Damien Allonsius, Daniel Song, Danielle Pintz, Danny Livshits, David Esiobu, Dhruv Choudhary, Dhruv Mahajan, Diego Garcia-Olano, Diego Perino, Dieuwke Hupkes, Egor Lakomkin, Ehab AlBadawy, Elina Lobanova, Emily Dinan, Eric Michael Smith, Filip Radenovic, Frank Zhang, Gabriel Synnaeve, Gabrielle Lee, Georgia Lewis Anderson, Graeme Nail, Gregoire Mialon,

Guan Pang, Guillem Cucurell, Hailey Nguyen, Hannah Korevaar, Hu Xu, Hugo Touvron, Iliyan Zarov, Imanol Arrieta Ibarra, Isabel Kloumann, Ishan Misra, Ivan Evtimov, Jade Copet, Jaewon Lee, Jan Geffert, Jana Vranes, Jason Park, Jay Mahadeokar, Jeet Shah, Jelmer van der Linde, Jennifer Billock, Jenny Hong, Jenya Lee, Jeremy Fu, Jianfeng Chi, Jianyu Huang, Jiawen Liu, Jie Wang, Jiecao Yu, Joanna Bitton, Joe Spisak, Jongsoo Park, Joseph Rocca, Joshua Johnstun, Joshua Saxe, Junteng Jia, Kalyan Vasuden Alwala, Kartikeya Upasani, Kate Plawiak, Ke Li, Kenneth Heafield, Kevin Stone, Khalid El-Arini, Krithika Iyer, Kshitiz Malik, Kuenley Chiu, Kunal Bhalla, Lauren Rantala-Yeary, Laurens van der Maaten, Lawrence Chen, Liang Tan, Liz Jenkins, Louis Martin, Lovish Madaan, Lubo Malo, Lukas Blecher, Lukas Landzaat, Luke de Oliveira, Madeline Muzzi, Mahesh Pasupuleti, Mannat Singh, Manohar Paluri, Marcin Kardas, Mathew Oldham, Mathieu Rita, Maya Pavlova, Melanie Kambadur, Mike Lewis, Min Si, Mitesh Kumar Singh, Mona Hassan, Naman Goyal, Narjes Torabi, Nikolay Bashlykov, Nikolay Bogoychev, Niladri Chatterji, Olivier Duchenne, Onur Çelebi, Patrick Alrassy, Pengchuan Zhang, Pengwei Li, Petar Vasic, Peter Weng, Prajjwal Bhargava, Pratik Dubal, Praveen Krishnan, Punit Singh Koura, Puxin Xu, Qing He, Qingxiao Dong, Ragavan Srinivasan, Raj Ganapathy, Ramon Calderer, Ricardo Silveira Cabral, Robert Stojnic, Roberta Raileanu, Rohit Girdhar, Rohit Patel, Romain Sauvestre, Ronnie Polidoro, Roshan Sumbaly, Ross Taylor, Ruan Silva, Rui Hou, Rui Wang, Saghar Hosseini, Sahana Chennabasappa, Sanjay Singh, Sean Bell, Seohyun Sonia Kim, Sergey Edunov, Shaoliang Nie, Sharan Narang, Sharath Raparthy, Sheng Shen, Shengye Wan, Shruti Bhosale, Shun Zhang, Simon Vandenhende, Soumya Batra, Spencer Whitman, Sten Sootla, Stephane Collot, Suchin Gururangan, Sydney Borodinsky, Tamar Herman, Tara Fowler, Tarek Sheasha, Thomas Georgiou, Thomas Scialom, Tobias Speckbacher, Todor Mihaylov, Tong Xiao, Ujjwal Karn, Vedanuj Goswami, Vibhor Gupta, Vignesh Ramanathan, Viktor Kerkez, Vincent Gonguet, Virginie Do, Vish Vogeti, Vladan Petrovic, Weiwei Chu, Wenhan Xiong, Wenyin Fu, Whitney Meers, Xavier Martinet, Xiaodong Wang, Xiaoqing Ellen Tan, Xinfeng Xie, Xuchao Jia, Xuewei Wang, Yaelle Goldschlag, Yashesh Gaur, Yasmine Babaei, Yi Wen, Yiwen Song, Yuchen Zhang, Yue Li, Yuning Mao, Zacharie Delpierre Coudert, Zheng Yan, Zhengxing Chen, Zoe Papakipos, Aaditya Singh, Aaron Grattafiori, Abha Jain, Adam Kelsey, Adam Shajnfeld, Adithya Gangidi, Adolfo Victoria, Ahuva Goldstand, Ajay Menon, Ajay Sharma, Alex Boesenberg, Alex Vaughan, Alexei Baevski, Allie Feinstein, Amanda Kallet, Amit Sangani, Anam Yunus, Andrei Lupu, Andres Alvarado, Andrew Caples, Andrew Gu, Andrew Ho, Andrew Poulton, Andrew Ryan, Ankit Ramchandani, Annie Franco, Aparajita Saraf, Arkabandhu Chowdhury, Ashley Gabriel, Ashwin Bharambe, Assaf Eisenman, Azadeh Yazdan, Beau James, Ben Maurer, Benjamin Leonhardi, Bernie Huang, Beth Loyd, Beto De Paola, Bhargavi Paranjape, Bing Liu, Bo Wu, Boyu Ni, Braden Hancock, Bram Wasti,

Brandon Spence, Brani Stojkovic, Brian Gamido, Britt Montalvo, Carl Parker, Carly Burton, Catalina Mejia, Changhan Wang, Changkyu Kim, Chao Zhou, Chester Hu, Ching-Hsiang Chu, Chris Cai, Chris Tindal, Christoph Feichtenhofer, Damon Civin, Dana Beaty, Daniel Kreymer, Daniel Li, Danny Wyatt, David Adkins, David Xu, Davide Testuggine, Delia David, Devi Parikh, Diana Liskovich, Didem Foss, Dingkang Wang, Duc Le, Dustin Holland, Edward Dowling, Eissa Jamil, Elaine Montgomery, Eleonora Presani, Emily Hahn, Emily Wood, Erik Brinkman, Esteban Arcaute, Evan Dunbar, Evan Smothers, Fei Sun, Felix Kreuk, Feng Tian, Firat Ozgenel, Francesco Caggioni, Francisco Guzmán, Frank Kanayet, Frank Seide, Gabriela Medina Florez, Gabriella Schwarz, Gada Badeer, Georgia Swee, Gil Halpern, Govind Thattai, Grant Herman, Grigory Sizov, Guangyi, Zhang, Guna Lakshminarayanan, Hamid Shojanazeri, Han Zou, Hannah Wang, Hanwen Zha, Haroun Habeeb, Harrison Rudolph, Helen Suk, Henry Aspegren, Hunter Goldman, Igor Molybog, Igor Tufanov, Irina-Elena Veliche, Itai Gat, Jake Weissman, James Geboski, James Kohli, Japhet Asher, Jean-Baptiste Gaya, Jeff Marcus, Jeff Tang, Jennifer Chan, Jenny Zhen, Jeremy Reizenstein, Jeremy Teboul, Jessica Zhong, Jian Jin, Jingyi Yang, Joe Cummings, Jon Carvill, Jon Shepard, Jonathan McPhie, Jonathan Torres, Josh Ginsburg, Junjie Wang, Kai Wu, Kam Hou U, Karan Saxena, Karthik Prasad, Kartikay Khandelwal, Katayoun Zand, Kathy Matosich, Kaushik Veeraraghavan, Kelly Michelena, Keqian Li, Kun Huang, Kunal Chawla, Kushal Lakhotia, Kyle Huang, Lailin Chen, Lakshya Garg, Lavender A, Leandro Silva, Lee Bell, Lei Zhang, Liangpeng Guo, Licheng Yu, Liron Moshkovich, Luca Wehrstedt, Madian Khabsa, Manav Avalani, Manish Bhatt, Maria Tsimpoukelli, Martynas Mankus, Matan Hasson, Matthew Lennie, Matthias Reso, Maxim Groshev, Maxim Naumov, Maya Lathi, Meghan Keneally, Michael L. Seltzer, Michal Valko, Michelle Restrepo, Mihir Patel, Mik Vyatskov, Mikayel Samvelyan, Mike Clark, Mike Macey, Mike Wang, Miquel Jubert Hermoso, Mo Metanat, Mohammad Rastegari, Munish Bansal, Nandhini Santhanam, Natascha Parks, Natasha White, Navyata Bawa, Nayan Singhal, Nick Egebo, Nicolas Usunier, Nikolay Pavlovich Laptev, Ning Dong, Ning Zhang, Norman Cheng, Oleg Chernoguz, Olivia Hart, Omkar Salpekar, Ozlem Kalinli, Parkin Kent, Parth Parekh, Paul Saab, Pavan Balaji, Pedro Rittner, Philip Bontrager, Pierre Roux, Piotr Dollar, Polina Zvyagina, Prashant Ratanchandani, Pritish Yuvraj, Qian Liang, Rachad Alao, Rachel Rodriguez, Rafi Ayub, Raghotham Murthy, Raghu Nayani, Rahul Mitra, Raymond Li, Rebekkah Hogan, Robin Battey, Rocky Wang, Rohan Maheswari, Russ Howes, Ruty Rinott, Sai Jayesh Bondu, Samyak Datta, Sara Chugh, Sara Hunt, Sargun Dhillon, Sasha Sidorov, Satadru Pan, Saurabh Verma, Seiji Yamamoto, Sharadh Ramaswamy, Shaun Lindsay, Shaun Lindsay, Sheng Feng, Shenghao Lin, Shengxin Cindy Zha, Shiva Shankar, Shuqiang Zhang, Shuqiang Zhang, Sinong Wang, Sneha Agarwal, Soji Sajuyigbe, Soumith Chintala, Stephanie Max, Stephen Chen, Steve Kehoe,

370

Steve Satterfield, Sudarshan Govindaprasad, Sumit Gupta, Sungmin Cho, Sunny Virk, Suraj Subramanian, Sy Choudhury, Sydney Goldman, Tal Remez, Tamar Glaser, Tamara Best, Thilo Kohler, Thomas Robinson, Tianhe Li, Tianjun Zhang, Tim Matthews, Timothy Chou, Tzook Shaked, Varun Vontimitta, Victoria Ajayi, Victoria Montanez, Vijai Mohan, Vinay Satish Kumar, Vishal Mangla, Vlad Ionescu, Vlad Poenaru, Vlad Tiberiu Mihailescu, Vladimir Ivanov, Wei Li, Wenchen Wang, Wenwen Jiang, Wes Bouaziz, Will Constable, Xiaocheng Tang, Xiaofang Wang, Xiaojian Wu, Xiaolan Wang, Xide Xia, Xilun Wu, Xinbo Gao, Yanjun Chen, Ye Hu, Ye Jia, Ye Qi, Yenda Li, Yilin Zhang, Ying Zhang, Yossi Adi, Youngjin Nam, Yu, Wang, Yuchen Hao, Yundi Qian, Yuzi He, Zach Rait, Zachary DeVito, Zef Rosnbrick, Zhaoduo Wen, Zhenyu Yang, and Zhiwei Zhao. 2024. The llama 3 herd of models. *CoRR*, abs/2407.21783.

Mor Geva, Roei Schuster, Jonathan Berant, and Omer Levy. 2021. Transformer feed-forward layers are key-value memories. In *Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing, EMNLP 2021, Virtual Event / Punta Cana, Dominican Republic, 7-11 November, 2021*, pages 5484–5495. Association for Computational Linguistics.

Soufiane Hayou, Nikhil Ghosh, and Bin Yu. 2024. Lora+: Efficient low rank adaptation of large models. *CoRR*, abs/2402.12354.

Neil Houlsby, Andrei Giurgiu, Stanislaw Jastrzebski, Bruna Morrone, Quentin de Laroussilhe, Andrea Gesmundo, Mona Attariyan, and Sylvain Gelly. 2019. Parameter-efficient transfer learning for NLP. In *Proceedings of the 36th International Conference on Machine Learning, ICML 2019, 9-15 June 2019, Long Beach, California, USA*, volume 97 of *Proceedings of Machine Learning Research*, pages 2790–2799. PMLR.

Edward J. Hu, Yelong Shen, Phillip Wallis, Zeyuan Allen-Zhu, Yuanzhi Li, Shean Wang, Lu Wang, and Weizhu Chen. 2022. Lora: Low-rank adaptation of large language models. In *The Tenth International Conference on Learning Representations, ICLR 2022, Virtual Event, April 25-29, 2022*. OpenReview.net.

Albert Q. Jiang, Alexandre Sablayrolles, Arthur Mensch, Chris Bamford, Devendra Singh Chaplot, Diego de Las Casas, Florian Bressand, Gianna Lengyel, Guillaume Lample, Lucile Saulnier, Lélio Renard Lavaud, Marie-Anne Lachaux, Pierre Stock, Teven Le Scao, Thibaut Lavril, Thomas Wang, Timothée Lacroix, and William El Sayed. 2023. Mistral 7b. *CoRR*, abs/2310.06825.

Damjan Kalajdzievski. 2023. A rank stabilization scaling factor for fine-tuning with lora. *CoRR*, abs/2312.03732.

Dawid Jan Kopiczko, Tijmen Blankevoort, and Yuki Markus Asano. 2023. Vera: Vector-based random matrix adaptation. *CoRR*, abs/2310.11454.

Xiang Lisa Li and Percy Liang. 2021. Prefix-tuning: Optimizing continuous prompts for generation. In *Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing, ACL/IJCNLP 2021, (Volume 1: Long Papers), Virtual Event, August 1-6, 2021*, pages 4582–4597. Association for Computational Linguistics.

Ilya Loshchilov and Frank Hutter. 2019. Decoupled weight decay regularization. In *7th International Conference on Learning Representations, ICLR 2019, New Orleans, LA, USA, May 6-9, 2019*. OpenReview.net.

Sadhika Malladi, Tianyu Gao, Eshaan Nichani, Alex Damian, Jason D. Lee, Danqi Chen, and Sanjeev Arora. 2023. Fine-tuning language models with just forward passes. In *Advances in Neural Information Processing Systems 36: Annual Conference on Neural Information Processing Systems 2023, NeurIPS 2023, New Orleans, LA, USA, December 10 - 16, 2023*.

Subhabrata Mukherjee, Arindam Mitra, Ganesh Jawahar, Sahaj Agarwal, Hamid Palangi, and Ahmed Awadallah. 2023. Orca: Progressive learning from complex explanation traces of GPT-4. *CoRR*, abs/2306.02707.

Colin Raffel, Noam Shazeer, Adam Roberts, Katherine Lee, Sharan Narang, Michael Matena, Yanqi Zhou, Wei Li, and Peter J. Liu. 2020. Exploring the limits of transfer learning with a unified text-to-text transformer. *J. Mach. Learn. Res.*, 21:140:1–140:67.

Victor Sanh, Albert Webson, Colin Raffel, Stephen Bach, Lintang Sutawika, Zaid Alyafeai, Antoine Chaffin, Arnaud Stiegler, Arun Raja, Manan Dey, M Saiful Bari, Canwen Xu, Urmish Thakker, Shanya Sharma Sharma, Eliza Szczechla, Taewoon Kim, Gunjan Chhablani, Nihal V. Nayak, Debajyoti Datta, Jonathan Chang, Mike Tian-Jian Jiang, Han Wang, Matteo Manica, Sheng Shen, Zheng Xin Yong, Harshit Pandey, Rachel Bawden, Thomas Wang, Trishala Neeraj, Jos Rozen, Abheesht Sharma, Andrea Santilli, Thibault Févry, Jason Alan Fries, Ryan Teehan, Teven Le Scao, Stella Biderman, Leo Gao, Thomas Wolf, and Alexander M. Rush. 2022. Multitask prompted training enables zero-shot task generalization. In *The Tenth International Conference on Learning Representations, ICLR 2022, Virtual Event, April 25-29, 2022*. OpenReview.net.

Vincenzo Scotti, Licia Sbattella, and Roberto Tedesco. 2024. A primer on seq2seq models for generative chatbots. *ACM Comput. Surv.*, 56(3):75:1–75:58.

Mojtaba Valipour, Mehdi Rezagholizadeh, Ivan Kobyzev, and Ali Ghodsi. 2023. Dylora: Parameter-efficient tuning of pre-trained models using dynamic search-free low-rank adaptation. In *Proceedings of the 17th Conference of the European Chapter of the Association for Computational Linguistics, EACL 2023, Dubrovnik, Croatia, May 2-6, 2023*, pages 3266–3279. Association for Computational Linguistics.

Thomas Wolf, Lysandre Debut, Victor Sanh, Julien Chaumond, Clement Delangue, Anthony Moi, Pierric Cistac, Tim Rault, Rémi Louf, Morgan Funtowicz, Joe Davison, Sam Shleifer, Patrick von Platen, Clara Ma, Yacine Jernite, Julien Plu, Canwen Xu, Teven Le Scao, Sylvain Gugger, Mariama Drame, Quentin Lhoest, and Alexander M. Rush. 2020. Transformers: State-of-the-art natural language processing. In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing: System Demonstrations, EMNLP 2020 - Demos, Online, November 16-20, 2020*, pages 38–45. Association for Computational Linguistics.

Qingru Zhang, Minshuo Chen, Alexander Bukharin, Pengcheng He, Yu Cheng, Weizhu Chen, and Tuo Zhao. 2023. Adaptive budget allocation for parameter-efficient fine-tuning. In *The Eleventh International Conference on Learning Representations, ICLR 2023, Kigali, Rwanda, May 1-5, 2023*. OpenReview.net.

Wayne Xin Zhao, Kun Zhou, Junyi Li, Tianyi Tang, Xiaolei Wang, Yupeng Hou, Yingqian Min, Beichen Zhang, Junjie Zhang, Zican Dong, Yifan Du, Chen Yang, Yushuo Chen, Zhipeng Chen, Jinhao Jiang, Ruiyang Ren, Yifan Li, Xinyu Tang, Zikang Liu, Peiyu Liu, Jian-Yun Nie, and Ji-Rong Wen. 2023. A survey of large language models. *CoRR*, abs/2303.18223.

## A  Source Code Availability

We share the source code associated with this paper for full reproducibility and transparency. All the source material to replicate the experiments is available on GitHub:

- L1RA: `https://github.com/raul-singh/L1RA/tree/dev-exp`;

- MEMORY-GELATO: `https://github.com/raul-singh/memory-gelato`.

## B  Evaluation setup

In this section, we provide the hyperparamers we used for in experimental evaluations to ensure full reproducibility. We report the hyperparameters we used with the OPENORCA data set in Table 4. In the table, we use the following notation:

- $r$ is the initial rank of L1RA, LoRA and ADALORA adapters;

- $\alpha$ is the scaling of L1RA, LoRA and ADALORA;

- $p_{dropout}$ is the dropout probability of L1RA, LoRA and ADALORA;

- $\eta$ is the learning rate;

- $\lambda$ is the regularisation coefficient or L1RA or ADALORA;

One important detail of our experiments is the choice of the optimiser, we implemented a variant of *AdamW* (Loshchilov and Hutter, 2019) (which is the most common optimiser for LLMs), to support decoupled regularisation for both $L_1$ and $L_2$ regularisations. We refer to this variant as *AdamE*, where the "E" refers to *ElasticNet*: the combined $L_1$ and $L_2$ regulariser[2]. The addition is the decoupled $L_1$ regularisation that avoids the update of the lasso constraint being scaled by the adaptive learning rate and momentum hyperparameters. This scaling affects negatively the shrinking of the parameters, showing it down.

Since we apply learning rate warm-up and cosine scheduling to shrink $\eta$ to zero, we find useful keep a separate constant learning rate for the parameters in the **c** vectors. To avoid introducing unnecessary hyperaparameters we use the same $\eta$ of the rest of the parameters, but whithout warm-up and scheduling.

[2]AdamE implementation `https://github.com/vincenzo-scotti/bitsandbytes/tree/dev-adame`

Table 4: OPENORCA hyperparameters.

| Model | Hyperparameter | Value |
|---|---|---|
| LLAMA 3.1 8B | Max. sequence length | 1024 tokens |
| | $r$ | 16 |
| | $\alpha$ | 16 |
| | $p_{dropout}$ | $10^{-1}$ |
| | Compute d-type | bfloat16 |
| | Attn. implementation | Flash attn. 2 (Dao, 2024) |
| | Optimiser | AdamE (paged, 32 bit) |
| | $\eta$ | $10^{-4}$ |
| | $\eta$ schedule | cosine |
| | $\eta$ warm-up ratio | 5% |
| | Max grad. norm | 1 |
| | Epochs | 1 |
| | Batch size | 4 |
| | Accum. steps | 4 |
| | $\lambda_{\text{L1RA}}$ | $10^{-3}$ |
| | $\eta_{\mathbf{c}}$ (L1RA) | $10^{-2}$ |
| | Rank update period (L1RA) | 5% training steps |
| | $\lambda_{\text{ADALORA}}$ | $10^{-3}$ |
| | $t_{init}$ (ADALORA) | 5% training steps |
| | $\Delta t$ (ADALORA) | 5% training steps |
| MISTRAL 7B V0.3 | Max. sequence length | 1024 tokens |
| | $r$ | 16 |
| | $\alpha$ | 16 |
| | $p_{dropout}$ | $10^{-1}$ |
| | Compute d-type | bfloat16 |
| | Attn. implementation | Flash attn. 2 (Dao, 2024) |
| | Optimiser | AdamE (paged, 32 bit) |
| | $\eta$ | $10^{-4}$ |
| | $\eta$ schedule | cosine |
| | $\eta$ warm-up ratio | 5% |
| | Max grad. norm | 1 |
| | Epochs | 1 |
| | Batch size | 4 |
| | Accum. steps | 4 |
| | $\lambda_{\text{L1RA}}$ | $10^{-3}$ |
| | $\eta_{\mathbf{c}}$ (L1RA) | $10^{-2}$ |
| | Rank update period (L1RA) | 5% training steps |
| | $\lambda_{\text{ADALORA}}$ | $10^{-3}$ |
| | $t_{init}$ (ADALORA) | 5% training steps |
| | $\Delta t$ (ADALORA) | 5% training steps |

We conducted all experiments on the same machine with the following hardware configuration:

- CPU: Intel Core i9-13900K;

- RAM: 64 GB;

- GPU: NVIDIA GeForce RTX 4090.

We used as much shared parameters across the three approaches we compare (L1RA, LoRA and ADALORA) as possible to have a fair comparison.

# Evaluating ASR in a Clinical Context : What Whisper Misses

**Haeeul Hwang[1,2], Eric Jordan[1], Deok-Hee Kim-Dufor[3], Christophe Lemey[3], Motasem Alrahabi[1]**

[1]Sorbonne Université, [2]Université de Paris Cité, [3]CHRU de Brest,

**Correspondence:** motasem.alrahabi@sorbonne-universite.fr

## Abstract

Automatic Speech Recognition (ASR) powered by AI is rapidly advancing and finding applications across a wide range of domains. However, its application within domain specific contexts still represents a challenge, with the specific issues varying from one context to another.

In this paper, we examine the application of OpenAI's Whisper ASR system in the context of psychiatric interviews. First, through analysis of error rates in the automatic transcriptions then through analysis of the most common errors present in the transcriptions we found that Whisper achieved a Word Error Rate of 0.25 but failed to transcribe filler words most often associated with patient's hesitations during speech. We find that systems such as Whisper show great promise for applications in clinical contexts . However, due to the importance of filler words and other filled pauses from a clinical perspective, its application should be accompanied with fine-tuning and verification by specialists to ensure the best outcomes.

## 1 Introduction

Despite improvements in ASR models, the latest models fail to achieve the same impressive out of the box results when applied to specific domains. This can result in increased errors for minority groups or those with speech disfluencies (Koenecke et al., 2020, 2024) or poor performance on languages that are under-represented within the training data of the aforementioned models (San et al., 2024).

Among these domains healthcare is a highly sensitive area where the accuracy of speech recognition is critical. In psychiatric consultations, subtle nuances in language can carry important clinical implications, and even small transcription errors may influence diagnosis, treatment, or documentation quality (Liebenthal et al., 2023). In these cases evaluating the models' true performance remains a challenge, since widely used evaluation metrics, such as Word Error Rate (WER) and Character Error Rate (CER) quantify surface-level inconsistencies between the recognized text and the reference transcript. However, these metrics do not account for semantic integrity, context relevance, or the clinical impact of recognition errors (Miner et al., 2020).

To address these challenges, this study analyses ASR outputs from psychiatric consultations. Specifically, we measure error rates across 3 different clinical groups and investigate the most common errors of the ASR models.

Our contributions are as follows:

1. We present a dataset of real-world psychiatric consultation transcripts and their ASR outputs (described further below).

2. We conduct an analysis of the recognition errors, with a focus on words that were omitted or deleted in the ASR-generated transcripts.

### 1.1 Prior work

The application of ASR in a clinical context remains underdeveloped, notably due to recognition errors that undermine reliability. Within the setting of psychotherapy consultations, Miner et al. (2020) introduced a three-pronged framework for assessing ASR, emphasizing that conventional metrics such as WER alone do not capture clinically relevant nuances. Their study compared human and ASR-generated transcriptions of therapy sessions, evaluating performance from three perspectives: general linguistic accuracy (WER and semantic distance), recognition of depression-related vocabulary, and accuracy on passages with harm-related language (e.g., self-harm or violence).

While the general performance analysis showed an average WER of 25% performance on harm-related speech was significantly lower with a 34% WER, indicating poor reliability in safety-critical scenarios. This underscores how ASR systems can provide transcriptions that are mostly accurate,

however they may miss finer details that are important for diagnosis. The authors conclude that these systems required further development before being ready for individual level safety surveillance.

Further work in the clinical domain has investigated the effectiveness of ASR systems for transcribing recordings of patients with Alzheimer's disease (Soroski et al., 2022; Akinrintoyo et al., 2025). Overall these studies showed promising results, with the transcripts produced being usable for distinguishing between patient and control groups (Soroski et al., 2022). However, the models' tendency to exclude filler words (e.g. *umm*, *uhh*) was noted, although this deficit could be made up for by fine-tuning Whisper on the patients' data (Akinrintoyo et al., 2025).

Additionally, some works have considered the use of ASR systems within the context of psychological experiments (Pfeifer et al., 2024; Ziman et al., 2018) finding error rates as low as 2.5 %, however these results only applied to studies with exclusively healthy participants (Pfeifer et al., 2024). Nonetheless, these works show the promise that these models can hold for application in psychological research.

## 2 Methodology

### 2.1 Dataset

| Patient Group | Female | Male | Total |
|---|---|---|---|
| AR | 22 | 19 | 41 |
| NAR | 5 | 5 | 11 |
| FEP | 3 | 4 | 7 |
| **Total** | **30** | **28** | **59** |

Table 1: Patient group counts by sex, including unknown and totals

| Group | Mean | Std | Median | Min | Max |
|---|---|---|---|---|---|
| Agg | 49 | 14 | 48 | 15 | 90 |
| NAR | 47 | 10 | 47 | 27 | 62 |
| AR | 50 | 14 | 48 | 16 | 90 |
| FEP | 48 | 18 | 53 | 15 | 74 |

Table 2: Overall and Group-Specific Recording Durations (in Minutes)

To assess ASR performance in psychiatric consultations, we compiled a dataset comprising audio

recordings of 59 (30 female, 28 male, 1 undisclosed) patients at ultra-high risk for psychosis (UHR) with a psychiatrist. All participants were native speakers of French. The interview was a semi-structured conversation with predetermined questions on each patient's problems such as their background, family, social relationships, socio-professional insertion, emotional interactions, complaints about symptoms, and other issues brought up by the patient.

The average recording duration was 49 minutes (see Table 2). Two trained assistants transcribed the entire utterances verbatim, including filled pauses, mispronunciations, and neologisms, following clear guidelines. An experienced linguist then reviewed the transcripts to correct spelling errors – such as homophones and accents – without altering the verbatim content. The participants were grouped into three clinical categories according to the Comprehensive Assessment of At-Risk Mental States (CAARMS) (Yung et al., 2005) : AR (At Risk for psychosis, 41 patients), NAR (Not At Risk, 11 patients), and FEP (First-Episode Psychosis, 7 patients), see Table 1.

### 2.2 Data Preprocessing

#### 2.2.1 Audio Processing

Prior to auto-transcription, all audio files – originally in formats such as WMA, M4A, MP4, and WAV – were converted to a uniform WAV format to ensure compatibility and consistency. The processed files were then transcribed using two ASR models: OpenAI's *Whisper medium* and *Whisper turbo* (Radford et al., 2023).

To improve evaluation accuracy, the start of each recording was trimmed to align the starting timepoints across files. Subsequently, a series of preprocessing steps were applied to both manual and ASR-generated transcripts before calculating error rates.

#### 2.2.2 Text Processing

1. **Special character removal** : We first removed non-verbal annotations and special characters (e.g., #, [, ]) that were used to distinguish between patient and clinician speech and to mark the proper names in the manual transcripts. This cleaned version of the reference transcripts was then used to compute baseline WER and CER scores.

2. **Normalization** : We applied the *Whisper nor-*

Figure 1: Comparison of WER Metrics Across Patient Groups for Whisper Medium and Whisper Turbo

| Processing Stage | Aggregate | NAR | AR | FEP | Kruskal–Wallis $p$-value |
|---|---|---|---|---|---|
| WER | 0.38 | 0.40 | 0.37 | 0.39 | 0.4889 |
| WER_normalized | 0.25 | 0.27 | 0.24 | 0.27 | 0.5981 |
| **WER_filtered** | **0.17** | **0.19** | **0.17** | **0.19** | 0.7969 |

Table 3: Mean WER Across Patient Groups With Kruskal–Wallis Test Results for Different Preprocessing Stages

*malizer* to both the reference and hypothesis transcripts in order to standardize formatting. This process involved removing punctuation, converting all text to lowercase, and normalizing common variants in transcription. The resulting metrics, calculated on these normalized texts, are reported as WER_normalized and CER_normalized.

3. **Stop word filtering** : Using the French language model from the SpaCy library, we removed stop words from both versions of the transcripts to emphasize semantically meaningful content.

4. **Single-letter token removal** : We excluded tokens consisting of a single character (e.g., "c" from *c'est*, "d" from *d'accord*), as these are often misrecognized and lack standalone semantic value.

5. **Filled pauses and filler word removal** : After analyzing omissions in the ASR-generated transcripts, we found that filled pauses and filler words were the most frequently omitted across all clinical groups (AR, NAR, and FEP). Based on this, we identified and removed 11 common words – such as *euh*, *bah*, *humm*, *oui*, *ok* and *non* – which consistently appeared missing in the ASR outputs. This enabled the calculation of refined metrics

(WER_filtered, CER_filtered) that better capture recognition quality in a clinical context.

For quantitative analysis between ASR-generated and manual transcripts, we used the jiwer library to compute standard evaluation metrics, including WER and CER. Beyond providing overall error rates, jiwer also facilitated the automatic extraction and categorization of specific error types such as substitutions, deletions, and insertions. This enabled a more fine-grained analysis of both the frequency and the nature of recognition errors across different models.

## 3 Results and Discussions

### 3.1 Model-Level WER Comparison

Figure 1 shows the distribution of WER scores for both Whisper models (medium and turbo), aggregated across all patients and preprocessing stages. Overall, the Whisper turbo model consistently outperformed the Whisper medium model across all metrics.

The most notable performance improvement occurred after removing filled pauses (e.g., *euh, bah, humm*), with WER decreasing by approximately 6–8 percentage points for both models. This highlights the strong influence of spontaneous speech markers on surface-level error rates in psychiatric dialogues.

## 3.2 Patient Group Comparison



Figure 2: Aggregate Distributions of WER by Patient Type

Figure 2 presents the WER distributions for each patient group: FEP, AR, and NAR. The differences between patient groups were more pronounced than those between models. Specifically, the FEP group exhibited the highest WERs throughout all preprocessing stages, indicating greater transcription difficulty.

## 3.3 Effect of Preprocessing

Initially, raw transcripts contained non-verbal annotations, special characters, and clinically irrelevant elements such as filled pauses and filler words, which can artificially inflate error rates if not accounted for. By systematically removing these components, we obtained a cleaner reference set that better reflected the meaningful content of the speech.

## 3.4 Statistical Analysis of Group Differences

To assess whether the differences in WER across clinical groups were statistically significant, we conducted a Kruskal–Wallis H test for each version of the metric: raw WER, WER_normalized, and WER_filtered. The results are summarized in Table ??.

None of the comparisons reached statistical significance, with p-values of 0.4889 (WER), 0.5981 (WER_normalized), and 0.7969 (WER_filtered), respectively. These results indicate that while median WER values differed slightly between the groups, the intra-group variability remained high, preventing clear group-level differentiation.

In particular, the wide spread of WER scores within the FEP group (as seen in Figure 2) suggests that individual differences – such as speech disfluency, cognitive state, and acoustic environment

– may play a stronger role than diagnostic group alone in shaping ASR performance.

As outlined above, our initial WER scores fell short of those presented in (Miner et al., 2020). This difference was largely explained by the presence of filler words within the reference transcripts that were ignored by Whisper. While these words can seem superfluous from a linguistic perspective, they can be deemed important from a clinical perspective. The frequency, manner and timing of these filled pauses can give an important insight into the mental state of a given patient. With this in mind, we plan to investigate fine-tuning transcription models to produce outputs that more closely resemble verbatim transcriptions, as discussed in (Akinrintoyo et al., 2025). Additionally, we aim to examine alternative evaluation metrics that may better capture clinically relevant transcription fidelity.

## 4 Conclusion

This study shows that the *medium* and *turbo* Whisper models perform well when applied in the context of clinical consultations in French, achieving an aggregate WER of 0.25 after applying Whisper normalisation (no statistical significance between patient groups was observed).

These results show that these models can be integrated into the interview process. Although the error rate would suggest that some human supervision and correction is still required, these automated transcriptions can provide a starting point that could significantly reduce the work load when transcribing interviews.

However, particular attention should still be paid to certain types of errors made by the ASR models. For example, we found 8 percentage points of the total error came from omitted filler words used during verbal pauses. These words can be essential to clinicians when diagnosing patients, and so these omissions are of the utmost importance. With this in mind, further work should follow the example of (Akinrintoyo et al., 2025) and investigate the possible fine tuning of ASR models to improve their accuracy vis-à-vis those linguistic details which are relevant to clinicians' diagnoses. This work provides valuable insight into ASR performance in a linguistic context such as French, as, to our knowledge, few studies have evaluated systems like Whisper in non-English clinical settings.

## Limitations

As outlined in Section 2.1 above, the dataset used for this experiment has a limited sample size and is not balanced in terms of the CAARMS groups. Table 1 shows that the number of AR patients is twice that of the other two groups combined. However, this distribution is typical of the clinical population, with a majority of patients falling into the at-risk category. The data presented here constitute a subset of our dataset for which both the recordings and transcriptions have been finalized. We are continuing to expand the dataset, with the aim of this experiment being to investigate the use of ASR to accelerate the transcription process.

Additionally, our analysis of transcription errors was conducted at the level of the entire transcription. This meant that no distinction was made between the patients' speech and the therapists' speech. Further work will aim to speaker turns to get a more accurate representation of performance. Finally, silent pauses were not analyzed in this study. In clinical interviews, the duration and timing of silences—alongside filled pauses—can provide meaningful cues about cognitive or emotional states. Incorporating silence duration into future ASR evaluations could offer a more comprehensive understanding of patient behavior and support finer-grained clinical insights.

## References

Emmanuel Akinrintoyo, Nadine Abdelhalim, and Nicole Salomons. 2025. WhisperD: Dementia Speech Recognition and Filler Word Detection with Whisper. *arXiv preprint*. ArXiv:2505.21551 [eess].

Allison Koenecke, Anna Seo Gyeong Choi, Katelyn X. Mei, Hilke Schellmann, and Mona Sloane. 2024. Careless whisper: Speech-to-text hallucination harms. In *Proceedings of the 2024 ACM Conference on Fairness, Accountability, and Transparency*, FAccT '24, page 1672–1681, New York, NY, USA. Association for Computing Machinery.

Allison Koenecke, Andrew Nam, Emily Lake, Joe Nudell, Minnie Quartey, Zion Mengesha, Connor Toups, John R. Rickford, Dan Jurafsky, and Sharad Goel. 2020. Racial disparities in automated speech recognition. *Proceedings of the National Academy of Sciences*, 117(14):7684–7689.

Einat Liebenthal, Michaela Ennis, Habiballah Rahimi-Eichi, Eric Lin, Yoonho Chung, and Justin T. Baker. 2023. Linguistic and non-linguistic markers of disorganization in psychotic illness. *Schizophrenia Research*, 259:111–120. Language and Speech Analysis in Schizophrenia and Related Psychoses.

Adam S. Miner, Albert Haque, Jason A. Fries, Scott L. Fleming, Denise E. Wilfley, G. Terence Wilson, Arnold Milstein, Dan Jurafsky, Bruce A. Arnow, W. Stewart Agras, Li Fei-Fei, and Nigam H. Shah. 2020. Assessing the accuracy of automatic speech recognition for psychotherapy. *npj Digital Medicine*, 3(1):82.

Valeria A. Pfeifer, Trish D. Chilton, Matthew D. Grilli, and Matthias R. Mehl. 2024. How ready is speech-to-text for psychological language research? Evaluating the validity of AI-generated English transcripts for analyzing free-spoken responses in younger and older adults. *Behavior Research Methods*, 56(7):7621–7631.

Alec Radford, Jong Wook Kim, Tao Xu, Greg Brockman, Christine Mcleavey, and Ilya Sutskever. 2023. Robust speech recognition via large-scale weak supervision. In *Proceedings of the 40th International Conference on Machine Learning*, volume 202 of *Proceedings of Machine Learning Research*, pages 28492–28518. PMLR.

Nay San, Georgios Paraskevopoulos, Aryaman Arora, Xiluo He, Prabhjot Kaur, Oliver Adams, and Dan Jurafsky. 2024. Predicting positive transfer for improved low-resource speech recognition using acoustic pseudo-tokens.

Thomas Soroski, Thiago da Cunha Vasco, Sally Newton-Mason, Saffrin Granby, Caitlin Lewis, Anuj Harisinghani, Matteo Rizzo, Cristina Conati, Gabriel Murray, Giuseppe Carenini, Thalia S Field, and Hyeju Jang. 2022. Evaluating web-based automatic transcription for alzheimer speech data: Transcript comparison and machine learning analysis. *JMIR Aging*, 5(3):e33460.

Alison R. Yung, Alison R. Yung, Hok Pan Yuen, Patrick D. Mcgorry, Lisa J. Phillips, Daniel Kelly, Margaret Dell'olio, Shona M. Francey, Elizabeth M. Cosgrave, Eoin Killackey, Carrie Stanford, Katherine Godfrey, and Joe Buckby. 2005. Mapping the onset of psychosis: The comprehensive assessment of at-risk mental states. *Australian & New Zealand Journal of Psychiatry*, 39(11-12):964–971. PMID: 16343296.

Kirsten Ziman, Andrew C. Heusser, Paxton C. Fitzpatrick, Campbell E. Field, and Jeremy R. Manning. 2018. Is automatic speech-to-text transcription ready for use in psychological experiments? *Behavior Research Methods*, 50(6):2597–2605.

# Tachelhiyt-Darija: a parallel speech corpus for two underrepresented languages

**Noureddine Atouf**
Chouaib Doukkali University
El Jadida, Morocco
`atouf.noureddine@ucd.ac.ma`

**Elsayed Issa**
Purdue University
West Lafayette, IN, USA
`esissa@purdue.edu`

**Said Ouzbayr**
Chouaib Doukkali University
El Jadida, Morocco
`ouzbayrsaid@gmail.com`

## Abstract

Despite recent advances in speech technology, several languages remain underrepresented. This linguistic disparity exacerbates the existing technological divide, resulting in limited access to speech-driven technologies. A key factor contributing to this challenge is the scarcity of datasets necessary to develop diverse speech recognition systems for low-sourced languages such as Amazigh and Moroccan Arabic. While the Amazigh language is emblematic of cultural identity and is deeply embedded in history, Darija remains the dialect spoken by the majority in Morocco. In this work, we introduce the first Tachelhiyt-Darija speech parallel corpus. A total of 24 Amazigh and 11 Moroccan Darija speakers recorded the parallel textual data, yielding a corpus of 2,772 audio segments. We also conducted benchmarking and fine-tuning of the Whisper ASR model. The results underscored the need for the development of datasets for under-resourced languages.

## 1 Introduction

According to the latest population census in Morocco (High Commission for Planning, 2024), Moroccan Arabic and Tamazight, both reckoned to be descendants of the Semitic branch of the Afro-Asiatic language family, are the most widely spoken languages in the country (Aissati et al., 2011). In a striking comparison, 92% of Moroccans use Darija (a term often interchangeably used with Moroccan Arabic), while only 25% of the population speaks Amazigh. At the sociolinguistic level, Moroccan Arabic and Amazigh are two dialectical forms spoken in different geographical areas and both have standardized forms: Modern Moroccan Arabic and Tamazight (Sadiqi, 2014; Youssi, 1992).

While Modern Standard Arabic (MSA) remains the country's first official language primarily utilized in formal speeches, administrative correspondences and documentation, news broadcast, and education as the medium of instruction, Moroccan Arabic and Tamazight have long maintained the status of the language of daily social interactions and conversations. In contrast to Darija, and particularly in 2001, Tamazight was recognized as a national heritage for all Moroccans, with the highest authority in the country (King Mohammed VI, 2001) issuing a decree to preserve, promote, and reinforce the Amazigh language and culture throughout the establishment of the Royal Institute of Amazigh Culture (IRCAM - Institut Royal de la Culture Amazighe). Amazigh was introduced into primary school curricula in 2003, marking a shift from its oral tradition to a formalized, codified language. This transition faced major challenges (Aissati et al., 2011; Ennaji, 2014), including standardizing the language across its three main varieties (namely, Tamazight, Tachelhiyt, and Tarifiyt) and selecting an appropriate script—ultimately, the Tifinagh script was adopted, evolving from the ancient Lybico-Berber alphabet.

The present work introduces the *Tachelhiyt-Darija* parallel speech corpus[1]. The corpus includes speech transcriptions in addition to speaker gender and dialect. In addition to describing our dataset, we developed baseline systems for automatic speech recognition (ASR). To summarize, our contributions are as follows. First, we introduce *Tachelhiyt-Darija*, a fully supervised speech dataset for two underrepresented dialects, labeled with transcriptions, dialect, and gender. Second, we evaluate Whisper as the state-of-the-art (SoTA) multilingual ASR model across the two dialects to assess its performance in recognizing underrepresented speech.

## 2 Related work

There is a limited body of research that develops speech data for underrepresented languages to

---

[1] `https://huggingface.co/datasets/NoureddineMOR/tachelhiyt-darija`

serve different purposes such as ASR models, machine translation, language learning, et cetera. Several scholars designed datasets for individual letters and digits for Amazigh spoken digit recognition tasks (Abakarim and Abenaou, 2023; Boulal et al., 2023; Hamidi et al., 2020; Telmem and Ghanou, 2018; El Ghazi et al., 2014; Satori and ElHaoussi, 2014).

On the word and sentence levels, El Ouahabi et al. (2017) developed a corpus for recognizing 520 spoken Amazigh words from 50 native Tarifiyt speakers, while (Oukas et al., 2024) created a database of Arabic vocal data by Tamazight speakers to train ASR models sensitive to Tamazight-accented Arabic. These efforts, along with (Daouad et al., 2023), focus on building word-based corpora to improve speech recognition systems and human-machine vocal interaction. Additionally, Mozilla Common Voice (Ardila et al., 2020) also contributes to this goal, offering a publicly available dataset with 398 Tamazight audio files for training and 159 for testing in version 17.0.

A different mode of database creation is attested in Moroccan Arabic. A variety of Moroccan Arabic databases have been developed, with the Darija Open Dataset (Outchakoucht and Es-Samaali, 2021) standing out for its 10,000 annotated entries featuring English translations and detailed linguistic information. Other studies, such as (Zaidani et al., 2024b) and (Zaidani et al., 2024a), constructed corpora through manual transcription and audio segmentation of YouTube content, while (Talafha et al., 2024) compiled 48 hours of transcribed data from North African and other Arabic dialect speakers—though Amazigh was notably excluded. Additional contributions include (Samih and Maier, 2016), (Ali et al., 2019), and (Labied et al., 2023), each aiming to support NLP, ASR, and speech-to-text translation tasks for Moroccan Arabic.

To our knowledge, there has not been an attempt to consider recording sentence stimuli in a parallel corpus consisting of two underrepresented languages, where the goal is to increase accessibility of the language in terms of the used script. In the context of the present experiment, we expand the Amazigh language speech corpus to include longer sequences instead of isolated words[2]. However, some publicly available datasets on Huggingface provide sentence recordings along with their written transcripts in Tifinagh[3] and phoentic symbols[45]. This study introduces a parallel corpus featuring two underrepresented languages—Tamazight and Darija—with vocal recordings transcribed in Arabic script for greater accessibility. Unlike existing repositories such as HuggingFace, which use Tifinagh or phonetic scripts for Tamazight, this work contributes a balanced, Arabic-script-based speech dataset for both languages.

## 3 Corpus Design

### 3.1 Parallel data and speech recoding

The process of designing the corpus was convoluted given the unavailability of written resources in Amazigh. Offline and online materials such as stories and written poetry exist in substantial amounts. Still, they are scripted in Tifinagh and require to be translated by a literate reader in the language. Moreover, visual content in Amazigh on the YouTube platform is abundant but is devoid of automatically generated subtitles. The absence of subtitled episodes with pure backgrounds determined our modus operandi with respect to data collection. The first stage of compiling the data was conducted through listing down numerous random (functional) sentences in Moroccan Arabic. The sentences were both generated by the authors and excerpted from the "Moroccan Arabic Textbook" designed by Peace Corps Morocco [6].

Two native speakers of Tachelhyit (the target Amazigh variety in the extant study) were, then, recruited to either literally translate the Moroccan Darija sentences into Tachelhyit. The rendered sentences were provided in the Arabic script. At the end, the data was further broken into sub-lists, which were sent to the participants who in turn recorded the stimuli in their first language.

The participants were two groups of Amazigh and Moroccan Arabic native speakers. The varying size of the groups, with the Amazigh participants forming the majority (n=24), reflected the study's sampling technique. Both male (n=14, 40%) and female (n=21, 60%) participants, whose age range differed across three major groups were included to

---

[2]https://huggingface.co/datasets/
Tamazight-NLP/tamawalt-n-imZZyann

[3]https://huggingface.co/datasets/
Tamazight-NLP/tosd

[4]https://huggingface.co/datasets/TifinLab/
moroccan_amazigh_asr

[5]https://huggingface.co/datasets/TifinLab/
tamazight_asr

[6]https://www.friendsofmorocco.org/Docs/Darija/
Moroccan%20Arabic%20textbook%202011.pdf

| Tachelhiyt | | | | Darija | | | |
|---|---|---|---|---|---|---|---|
| **audio** | **transcript** | **speaker** | **gender** | **audio** | **transcript** | **speaker** | **gender** |
| t_1.wav | أزول فلاون | t_sp14 | male | d_1.wav | السلام عليكم | d_sp14 | male |
| t_2.wav | رادفتوخ ناف أوزكا | t_sp14 | male | d_2.wav | غادي نمشي بعد غدا | d_sp14 | male |
| t_3.wav | لوقت ن الطوبيس | t_sp12 | female | d_3.wav | توقيت الطوبيسات | d_sp12 | female |

Table 1: Sample of the Tachelhiyt-Darija speech parallel corpus

diversify the spoken form of our data. The majority of participants (57.10%) were under 19 years old, followed by 25.70% aged 20–29, and 17.10% aged 30–40. Table 2 shows the count and percentage rates for such demographic information as gender and participants' native tongue. For the age factor, the table displays the mean, range and standard deviation:

| Variables | Values | N | % | M | R | SD |
|---|---|---|---|---|---|---|
| Gender | Male | 14 | 40 | | | |
| | Female | 21 | 60 | | | |
| Native L. | Amazigh | 24 | 68 | | | |
| | Moroccan | 11 | 31 | | | |
| Age | | | | 22 | 19 | 6 |
| Education | high school | 21 | 60 | | | |
| | undergrad. | 5 | 14 | | | |
| | postgraduate | 9 | 25 | | | |

Table 2: Demographic Variables

## 3.2 Data preparation and segmentation

All participants recorded the data in the wild without any professional equipment, resulting in audio files with different sampling rates, 44.100 and 48.000, for Tachelhiyt and Darija respectively. The participants were asked to make a five-second pause between a recorded string and the following one. This allowed for optimal splitting of the data using *AudioSegment*[7] to extract the speech segments. The resulted data was verified by the experimenters to make sure the audio files were segmented properly and were aligned well with the corresponding transcripts.

Each audio segment was converted to a single channel 16 kHz 16-bit PCM encoded WAV files using the FFmpeg library (Tomar, 2006). For each dialect, there is a comma-separated (CSV) file con-

---

[7] https://github.com/jiaaro/pydub

taining two columns with the audio file name as the first column and the transcript as the second column. The metadata file includes eight columns; four for each dialect as shown in Table 1. The metadata further facilitated the filtering of specific speakers on the basis of their number and gender.

A total of 2772 speech segments made up the developed parallel corpus in Tachelhiyt and Draija.The duration of the Tachelhiyt recordings is 71.68 minutes, with an average audio segment length of 3.11 seconds, while the duration of the Darija recordings is 61.84 minutes, with an average segment length of 2.68 seconds.

## 4 Experiments

We evaluated Whisper (Radford et al., 2023) performance using zero-shot and full finetuning evaluation. Our primary goal is to report the importance of speech datasets by benchmarking one of the state-of-the-art ASR models.

### 4.1 Zero-shot evaluation

We evaluated Whisper small (242M) and large-v3 (1.54B) in a zero-shot setting using the test dataset (i.e., 300 speech samples). Then, we reported Word Error Rate (WER) and Character Error Rate (CER).

| Model | Tachelhiyt | Darija |
|---|---|---|
| | wer/cer | wer/cer |
| Baseline (small) | 127.09/80.85 | 89.72/60.47 |
| Baseline (large) | 145.54/85.86 | 76.08/34.44 |

Table 3: Results of Word Error Rate (WER) and Character Error Rate (CER) for the baseline models.

Table 3 shows the performance results for the baseline models on Tachelhiyt and Darija. The results reveal notable discrepancies across both languages and model sizes. For Tachelhiyt, the WER/CER increased from 127.09/80.85 in the

small model to 145.54/85.86 in the large model, indicating a performance degradation with the larger model. In contrast, Darija showed improvement with model scaling: the WER/CER decreased from 89.72/60.47 in the small model to 76.08/34.44 in the large model. These results suggest that the large model is more effective for Darija but less suited for Tachelhiyt. This is possibly due to differences in linguistic structure, training data distribution, or model overfitting.

## 4.2 Full Finetuning

For fine-tuning, we sampled audio data at a sampling rate of 16 kHz. All experiments were conducted on a single-node Google Colab instance equipped with an A100 GPU. We fine-tuned two Whisper small models separately for Tachelhiyt and Darija, using identical hyperparameters for both models. The fine-tuning process employed a training batch size of 16 and an evaluation batch size of 8. We set the learning rate to 1e-5 and trained for a maximum of 1000 steps, equivalent to approximately 15 epochs. These hyperparameters were chosen to accommodate the relatively small size of the available training data. For decoding, we used a maximum sequence length of 225 tokens. No additional post-processing steps were applied to the decoded outputs. Although we evaluated both the Whisper small and large-v3 models in a zero-shot setting, we opted to fine-tune the small model due to the limited size of the available dataset.

| Model | Tachelhiyt | Darija |
|---|---|---|
| | wer/cer | wer/cer |
| (small) | 7.45/3.25 | 4.26/1.38 |

Table 4: Results of Word Error Rate (WER) and Character Error Rate (CER) for the finetuned models.

As Table 4 shows, the fine-tuned small model demonstrates a substantial performance improvement over the baseline, achieving significantly lower error rates across both languages. For Tachelhiyt, WER and CER dropped to 7.45 and 3.25, respectively, while Darija saw even lower error rates of 4.26 (WER) and 1.38 (CER), indicating the effectiveness of fine-tuning in enhancing model accuracy for both language varieties.

The findings have important implications for the development and use of ASR systems in underrepresented languages. First, the baseline performance highlights the challenges that pre-trained ASR models face when deployed in low-resource languages such as Tachelhiyt and Darija. The very high WER and CER across board consistently show that such models are not inherently able to generalize to linguistically diverse contexts without some form of adaptation. This finding further supports the need for the creation of datasets and finetuning to broaden the applicability of ASR technologies to resource-poor languages, which mitigates the existing technological gap.

## 5 Conclusion

To conclude, this study introduces the first sequential parallel corpus of Tachelhiyt-Darija, comprising 71.68 minutes of Tachelhiyt and 61.84 minutes of Darija. Based on our review of the literature, this is the first parallel dataset of North African speech data from two less represented dialects. It has the potential to be used in Automatic Speech Recognition (ASR), speech-to-speech translation, and machine translation, with further applications for cross-lingual studies between these two languages. We have also applied and validated the Whisper ASR model by means of the utilized dataset for benchmarking and fine-tuning. The findings reveal the importance of creating language-specific datasets for less-documented languages in order to improve the current state of the art in speech technologies as well as enhance linguistic accessibility.

## Limitations

This study has certain limitations that should be acknowledged. The first limitation is that the quality and consistency of the audio data may have been affected by variations in participants' recording equipment, eventually leading to noise and possible background interference. This is ascribed to the small size and parallel nature of the dataset, which may have limited the scope of training and evaluation for the fine-tuned models, consequently limiting the generalizability of the results. However, the results are quite insightful. A larger and more diverse corpus could very well enhance the performance and robustness of the system.

In this vein, future work will involve the release of an expanded dataset— another parallel corpus—aimed at increasing both the quantity and quality of data for Tamazight three varieties and Darija, thereby enhancing ASR technology for these linguistic groups. In addition, future research

will bring structural and linguistic differences (i.e., phonological and morphological components) that characterize these two languages. This should further improve the performance of the ASR system.

## Ethics Statement

In developing Tachelhiyt-Darija corpus, we adhered to ethical principles to ensure responsible and respectful use of data. All speech data used in this study were collected in strict accordance with ethical research standards. Participants were fully informed about the purpose of the study, the nature of the data being collected, and their right to withdraw at any time without consequence. No personally identifiable information was collected, and all audio recordings were anonymized to protect participant privacy. The collected data are used exclusively for academic and research purposes related to language resource development and are stored securely to prevent unauthorized access.

## References

Fadwa Abakarim and Abdenbi Abenaou. 2023. Enhancing amazigh speech recognition system with mfdwc-svm. In *International Conference on Computational Science and Its Applications*, pages 471–488. Springer.

Abdelilah El Aissati, Susy Karsmakers, and Jeanne Kurvers. 2011. "we are all beginners": Amazigh in language policy and educational practice in morocco. *Compare: A Journal of Comparative and International Education*, 41(2):211–227.

Ahmed Ali, Suwon Shon, Younes Samih, Hamdy Mubarak, Ahmed Abdelali, James Glass, Steve Renals, and Khalid Choukri. 2019. The mgb-5 challenge: Recognition and dialect identification of dialectal arabic speech. In *2019 IEEE Automatic Speech Recognition and Understanding Workshop (ASRU)*, pages 1026–1033. IEEE.

Rosana Ardila, Megan Branson, Kelly Davis, Michael Kohler, Josh Meyer, Michael Henretty, Reuben Morais, Lindsay Saunders, Francis Tyers, and Gregor Weber. 2020. Common voice: A massively-multilingual speech corpus. In *Proceedings of the Twelfth Language Resources and Evaluation Conference*, pages 4218–4222, Marseille, France. European Language Resources Association.

Hossam Boulal, Mohamed Hamidi, Mustapha Abarkan, and Jamal Barkani. 2023. Amazigh spoken digit recognition using a deep learning approach based on mfcc. *International journal of electrical and computer engineering systems*, 14(7):791–798.

Mohamed Daouad, Fadoua Ataa Allah, and El Wardani Dadi. 2023. An automatic speech recognition system for isolated amazigh word using 1d & 2d cnn-lstm architecture. *International Journal of Speech Technology*, 26(3):775–787.

Ahmed El Ghazi, Cherki Daoui, and Najlae Idrissi. 2014. Automatic speech recognition for tamazight enchained digits. *World Journal Control Science and Engineering*, 2(1):1–5.

Safâa El Ouahabi, Mohamed Atounti, and Mohamed Bellouki. 2017. A database for amazigh speech recognition research: Amzsrd. In *2017 3rd International Conference of Cloud Computing Technologies and Applications (CloudTech)*, pages 1–5. IEEE.

Moha Ennaji. 2014. *Multiculturalism and Democracy in North Africa: Aftermath of the Arab Spring*. Routledge.

Mohamed Hamidi, Hassan Satori, Ouissam Zealouk, and Khalid Satori. 2020. Amazigh digits through interactive speech recognition system in noisy environment. *International Journal of Speech Technology*, 23(1):101–109.

High Commission for Planning. 2024. Morocco population and housing census.

Maria Labied, Abdessamad Belangour, and Mouad Banane. 2023. Darija-c: towards a moroccan darija speech recognition and speech-to-text translation corpus. In *2023 1st International Conference on Advanced Innovations in Smart Cities (ICAISC)*, pages 1–4. IEEE.

Nourredine Oukas, Tiziri Chabi, and Tilelli Sari. 2024. A novel dataset for arabic speech recognition recorded by tamazight speakers. *Authorea Preprints*.

Aissam Outchakoucht and Hamza Es-Samaali. 2021. Moroccan dialect-darija-open dataset. *arXiv preprint arXiv:2103.09687*.

Alec Radford, Jong Wook Kim, Tao Xu, Greg Brockman, Christine McLeavey, and Ilya Sutskever. 2023. Robust speech recognition via large-scale weak supervision. In *International conference on machine learning*, pages 28492–28518. PMLR.

Fatima Sadiqi. 2014. Berber and language politics in the moroccan educational system. In Moha Ennaji, editor, *Multiculturalism and Democracy in North Africa: Aftermath of the Arab Spring*. Routledge.

Younes Samih and Wolfgang Maier. 2016. An arabic-moroccan darija code-switched corpus. In *Proceedings of the Tenth International Conference on Language Resources and Evaluation (LREC'16)*, pages 4170–4175.

Hassan Satori and Fatima ElHaoussi. 2014. Investigation amazigh speech recognition using cmu tools. *International Journal of Speech Technology*, 17:235–243.

Bashar Talafha, Karima Kadaoui, Samar Mohamed Magdy, Mariem Habiboullah, Chafei Mohamed Chafei, Ahmed Oumar El-Shangiti, Hiba Zayed, Rahaf Alhamouri, Rwaa Assi, Aisha Alraeesi, et al. 2024. Casablanca: Data and models for multidialectal arabic speech recognition. *arXiv preprint arXiv:2410.04527*.

Meryam Telmem and Youssef Ghanou. 2018. Estimation of the optimal hmm parameters for amazigh speech recognition system using cmu-sphinx. *Procedia Computer Science*, 127:92–101.

Suramya Tomar. 2006. Converting video formats with ffmpeg. *Linux Journal*, 2006(146):10.

Abderrahim Youssi. 1992. *Grammaire et lexique de l'arabe marocain moderne*. Wallada, Casablanca.

Hajar Zaidani, Abderrahim Maizate, Mohammed Ouzzif, and Rim Koulali. 2024a. Building a corpus for the underexplored moroccan dialect (cfmd) through audio segmentations. *Revue d'Intelligence Artificielle*, 38(3):857.

Hajar Zaidani, Abderrahim Maizate, Mohammed Ouzzif, and Rim Koulali. 2024b. Cfmd: Corpus for moroccan dialect as under researched dialect. In *Future of Information and Communication Conference*, pages 61–69. Springer.

# From Outliers to Topics in Language Models: Anticipating Trends in News Corpora

**Evangelia Zve, Benjamin Icard, Alice Breton,**
**Lila Sainero, Gauvain Bourgne, and Jean-Gabriel Ganascia**

LIP6, Sorbonne University, CNRS, France

## Abstract

This paper examines how outliers, often dismissed as noise in topic modeling, can act as weak signals of emerging topics in dynamic news corpora. Using vector embeddings from state-of-the-art language models and a cumulative clustering approach, we track their evolution over time in French and English news datasets focused on corporate social responsibility and climate change. The results reveal a consistent pattern: outliers tend to evolve into coherent topics over time across both models and languages.

## 1 Introduction

As information ecosystems become increasingly dynamic, the early identification of emerging trends in news media remains a key challenge for natural language processing. Topic modeling, which clusters semantically similar documents to uncover latent themes, plays a central role in this task. Early approaches, most notably Latent Dirichlet Allocation (LDA) (Blei et al., 2003), introduced a probabilistic framework to infer latent topics from textual documents (Hoyle et al., 2022). More recent embedding-based methods, such as BERTopic (Grootendorst, 2022), represent documents as dense vector embeddings, enabling more contextualized representations and yielding more coherent topics in dynamic corpora such as online news content (Babalola et al., 2024).

Unlike partition-based clustering methods often used for clustering vector embeddings, such as KMeans (Hartigan and Wong, 1979), or probabilistic topic models like LDA, both of which assign every document to a topic, HDBSCAN (Campello et al., 2015) is a density-based clustering algorithm that explicitly labels low-density points as *outliers*. These documents, which do not fit into any topical cluster, are often treated as noise and excluded from downstream analysis.

Challenging the assumption that outliers are mere noise, we explore the hypothesis that outliers, documents not assigned to any cluster, may serve as early signals of emerging topics. We employ a cumulative clustering approach using BERTopic with HDBSCAN, tracing how isolated documents evolve and whether they are gradually integrated into clusters as their narratives gain salience. To aid interpretability, we also analyze lexical and stylistic features of outliers and their role in cluster integration.

To conduct our analysis, we use two news corpora. The first, in French, is a manually curated dataset documenting a corporate social responsibility dispute which serves as a pilot study. The second, in English, focuses on climate change and is used for replication. Both corpora are topically constrained, span continuous time periods, and provide full-text coverage, allowing to control for topical and timeline gaps.

Section 2 reviews related work. Section 3 details the full experimental setting, with a particular focus on the methodology. Section 4 presents the French study and results on outlier conversion. Section 5 reports replication results in English. Findings in both languages are discussed and compared in Section 6. Section 7 concludes and outlines future directions.

## 2 Related Work

Topic modeling is widely applied across various domains, including corporate social responsibility (Lee et al., 2023) and climate change (Ylä-Anttila et al., 2022), in both traditional and social media contexts (Laureate et al., 2023). The field's methodological evolution, from probabilistic approaches like LDA (Blei et al., 2003) to embedding-based methods such as BERTopic (Grootendorst, 2022), has improved semantic coherence. However, while outliers have been often treated as noise (Alattar and Shaalan, 2021), their role in sig-

naling emerging topics remains an underexplored area of research.

Research in temporal topic analysis has evolved from early techniques like burst detection (Chen et al., 2016) and term-frequency-based change point identification (Yao et al., 2021) to more recent approaches tracking semantic drift (Jung et al., 2020) and transformer-based dynamic modeling (Karakkaparambil et al., 2024; Boutaleb et al., 2024). While these methods effectively capture shifts in established topics, they typically overlook sparse outliers, documents that may precede and predict emerging themes before they coalesce into detectable clusters.

This relates to clustering methodology. While probabilistic topic models like LDA assign soft cluster memberships, and partition-based algorithms such as KMeans (Hartigan and Wong, 1979) enforce hard assignments, both approaches assume that every document belongs to a cluster. In contrast, density-based methods like HDBSCAN (Campello et al., 2015) and OPTICS (Ankerst et al., 1999) explicitly identify outliers as low-density points that do not belong to any cluster. Unlike general anomaly detection techniques (e.g., Isolation Forest (Liu et al., 2008), Local Outlier Factor (Breunig et al., 2000)), which detect outliers without considering the topical coherence of thematically structured corpora, HDBSCAN's built-in outlier detection aligns more closely with semantic structure. This allows to track how semantically isolated documents may evolve into coherent topic clusters over time.

This paper examines whether outliers can serve as early signals of emerging topics. By tracking their integration into clusters over time via cumulative clustering, we aim to complement existing work focused on stable topic structures.

## 3 Experimental Setting

### 3.1 Hypothesis

While topic modeling and document clustering have been extensively studied, the role of outliers in the dynamic formation of topics has not yet been explored. To address this gap, we propose the following hypothesis:

$\mathcal{H}$: *In topic-based cumulative clustering of news articles, topics emerge or are reinforced in part through the assimilation of outliers—that is, documents initially unclustered that later become part of coherent topic clusters.*

This hypothesis assumes that topic formation in cumulative clustering reflects a gradual process of semantic integration, in which outliers may act as early signals of emerging or evolving topics.

### 3.2 Models

To test $\mathcal{H}$, we use nine open-source embedding models with diverse transformer architectures and language capabilities. Model selection was guided by performance on the Massive Text Embedding Benchmark (MTEB) (Muennighoff et al., 2022), as reported on the Hugging Face leaderboard[1] as of September 16, 2024. Table 3 (Appendix A.2) summarizes the selected models.

### 3.3 Methodology

The methodology involves four main steps. First, we project news articles into a semantic space using language model embeddings. We then apply dimensionality reduction to enable efficient clustering and address the *curse of dimensionality* (Köppen, 2000). Subsequently, we perform cumulative clustering over 20 monthly time windows and evaluate clustering quality to determine the optimal experimental configuration. Based on this setup, we test $\mathcal{H}$ concerning outlier-to-topic conversion and assess its robustness through inter-model agreement. Finally, we analyze lexical and stylistic features to interpret differences between converted and non-converted outliers.

#### 3.3.1 Data Preparation

Each news article is represented using dense vector embeddings generated from nine pre-trained language models. For each document, we compute embeddings from three variants: body text, headline, and full article (both headline and body text). This projects articles into a high-dimensional semantic space, where distances reflect semantic similarity. We apply Uniform Manifold Approximation and Projection (UMAP) (McInnes et al., 2018) to reduce the dimensionality of embeddings prior to clustering. Output dimensions are varied across 2D, 3D, 5D, and 10D. UMAP is chosen over Principal Component Analysis (PCA) (Wold, 1987) due to its ability to preserve both local and global structure, which is important for identifying fine-grained topic distinctions and local outliers (Atzberger et al., 2023).

---

[1] https://huggingface.co/spaces/mteb/leaderboard

### 3.3.2 Cumulative Clustering

We employ cumulative clustering (iterative topic modeling over expanding time windows) across 20 monthly intervals. At each step, documents from the current and all prior months are clustered jointly using BERTopic with HDBSCAN (McInnes et al., 2017). This density-based algorithm assigns documents to clusters or labels them as outliers via the GLOSH algorithm (Campello et al., 2015), which identifies low-density regions by comparing a point's local density to its neighbors. Documents labeled -1 are classified as outliers and excluded from clusters. To test $\mathcal{H}$, we track whether these outliers transition to inliers (i.e., join a cluster) in subsequent windows, thereby signaling emergent topics.

The clustering quality is evaluated using the silhouette score (Shahapure and Nicholas, 2020), which measures cluster cohesion and separation. Scores above 0.7 are considered strong, 0.5–0.7 moderate, and below 0.25 weak. To evaluate clustering over time, we compute the mean and median silhouette scores across all time windows, and then aggregate these globally across all models. We compare the nine selected embedding models, content variants (headline, body, full article), and UMAP settings to ensure robustness. Based on these comparisons, we select the configuration with the highest silhouette score and proceed with our methodology to verify our hypothesis.

### 3.3.3 Outlier-to-Topic Conversion

Under hypothesis $\mathcal{H}$, we evaluate whether outliers contribute to the formation of new topics or the reinforcement of existing ones. We compute, for each model, the proportion of outliers that later become topic inliers, and assess robustness via the rescaling method of Icard et al. (2024), which measures whether $\mathcal{H}$ is consistently validated for the *same* outliers across models. Specifically, for articles identified as outliers by *all* models (at some point in their time window), we compute the proportion $x$ of models that validate $\mathcal{H}$ and rescale it as $a = |2x - 1|$. This transformation captures consensus independently of polarity (as in Cohen's kappa): both $x = 1$ (unanimous validation) and $x = 0$ (unanimous rejection) yield maximal agreement $a = 1$, while $x = 0.5$ corresponds to minimal agreement $a = 0$, since models are evenly split in this case.

### 3.3.4 Lexicon and Writing Style Analysis

As an attempt to explain the conversions observed, we first controlled for potential topical differences between converted and non-converted outliers using word-level TfidfVectorizer scores (Qaiser and Ali, 2018), hereafter referred to as TF-IDF. Let $w$ be a word and let $\text{TFIDF}_g(w)$ denote its average TF-IDF score in group $g \in \{\mathcal{H}, \text{not } \mathcal{H}\}$, where $\mathcal{H}$ corresponds to outliers that were integrated into topic clusters ("converted"), and not $\mathcal{H}$ to those that remained isolated ("non-converted"). To capture the differential lexical salience of word $w$ across the two groups, we define the delta TF-IDF as:

$$\Delta \text{TFIDF}(w) = \text{TFIDF}_{\mathcal{H}}(w) - \text{TFIDF}_{\text{not } \mathcal{H}}(w) \quad (1)$$

In addition, we investigated variation beyond lexical content by analyzing stylistic differences between converted and non-converted outliers using the stylometric framework introduced by Terreau et al. (2021), which quantifies eight core stylistic dimensions. These include the relative frequency of *function words* (e.g., prepositions, conjunctions, auxiliaries), *punctuation marks* (e.g., periods, commas), *numbers*, and *named entities* (e.g., persons, organizations) per sentence; distributions of *part-of-speech tags* (e.g., nouns, verbs, adjectives); and averages of *structural features* (e.g., word length, word frequency, syllables per word). The framework also incorporates *lexical complexity metrics* (e.g., Yule's K (Yule, 2014), Shannon entropy (Shannon, 1948)) and *readability indices* (e.g., Flesch-Kincaid Grade Level (Kincaid et al., 1975)).

## 4 Pilot Study

### 4.1 French Dataset

We constructed a dataset for the pilot study, referred to as TP, consisting of 102 French news articles that we manually collected and curated. The articles document a controversy involving the major energy company *TotalEnergies* and the prestigious French Grande École *École Polytechnique*, who planned to build a research center on the university's Saclay campus. The project drew both support, citing its contribution to energy research, and criticism, focused on academic independence and environmental impact. The TP dataset covers the full timeline of media coverage, from December 2018 to August 2024, and includes documents from official sources, mainstream media, partisan

outlets, opinion sections, and NGOs. It captures the entire development of the story, without topical or temporal gaps.

## 4.2 Topic-Based Clustering

We applied topic-based clustering to the TP dataset using the methodology described in subsection 3.3. Figure 1 presents the cumulative clustering output generated by the `Solon-embeddings-large-0.1` model. The figure shows a 2D representation derived from 10D `UMAP` projections of document embeddings across nine time windows, illustrating topic structure and outlier transitions over time.



Figure 1: 2D Scatter plot of the cumulative clustering obtained on TP (after `UMAP` 10D reduction) over nine time windows, using `Solon-embeddings-large-0.1`. Outliers are indicated with black $\times$ and topics in blue and green.

Across all nine models and `UMAP` dimensions, clustering quality is consistent, with mean and median silhouette scores above 0.5 (range: -1 to 1). On average, body-text embeddings yield higher-quality clusters than headline or full-article representations. `UMAP` with 10 dimensions outperforms the 2D, 3D, and 5D settings. Among models, `Solon-embeddings-large-0.1` achieves the highest scores, while `xlm-roberta-large` performs the worst. Based on these findings, we evaluate Hypotheses $\mathcal{H}$ on TP using `UMAP-10D` and body-text embeddings.

## 4.3 Outlier Behavior

To evaluate Hypothesis $\mathcal{H}$, we computed, for each model, the proportion of outliers that later became inliers during cumulative clustering. Figure 2 shows the mean validation score per model.

Table 1: Mean silhouette scores per model for `UMAP` 10D using the body text of the TP dataset. Bold values indicate the models achieving the best silhouette score for each document type. (See full results in A.3.1.)

| Model | UMAP 10D | | |
|---|---|---|---|
| | Headline | Body | Full Article |
| `multilingual-e5-large` | **0.6065** | 0.5519 | 0.5689 |
| `e5-base-v2` | 0.5592 | 0.5350 | 0.4846 |
| `sentence-camembert-base` | 0.5990 | 0.5850 | 0.6167 |
| `all-MiniLM-L12-v2` | 0.5654 | 0.5846 | 0.5349 |
| `Solon-embeddings-large-0.1` | 0.5772 | 0.6694 | 0.5553 |
| `xlm-roberta-large` | 0.4941 | 0.4802 | 0.4424 |
| `all-roberta-large-v1` | 0.5525 | 0.6258 | 0.5759 |
| `multilingual-mpnet-base-v2` | 0.5391 | 0.5923 | 0.6865 |
| `distilbert-base-uncased` | 0.3670 | **0.9373** | **0.8895** |
| Mean | 0.5400 | **0.6180** | 0.5993 |
| Median | 0.5417 | **0.6183** | 0.5756 |



Figure 2: Mean number of outliers per model that validate prediction $\mathcal{H}$ on TP by converting into topic inlier at some time point (specific to each model). Each colored bar represents the mean of each model.

The average validation score of $\mathcal{H}$ across models on TP is high, with a mean of 0.80. As expected, models trained or fine-tuned on French perform strongly: `Solon-embeddings-large-0.1` achieves perfect validation (1.0), and `sentence-camembert-base` scores 0.92. Among English-language models, `e5-base-v2` shows intermediate performance (0.68), while several others yield unexpectedly strong results: `all-MiniLM-L12-v2` (0.74), `all-roberta-large-v1` (0.84), and `distilbert-base-uncased` (1.0). Multilingual models show mixed performance: `xlm-roberta-large` scores moderately (0.65), whereas `paraphrase-multilingual-mpnet-base-v2` (0.82) and `multilingual-e5-large` (0.78) achieve high scores. Overall, model-level validation of $\mathcal{H}$ ranges from moderate to perfect, with a relatively uniform distribution.

Across models, outlier-to-inlier conversion rates are highest in early clustering phases (64.58%–100% in late 2020), followed by a ta-

pering trend with persistent outliers in later periods. As detailed in Appendix A.3.2, some models exhibit stable integration while others decline over time. Despite these intra-model fluctuations, the general pattern of early integration supports $\mathcal{H}$ across temporal windows.

To test whether this holds beyond model variation, we computed inter-agreement using the rescaling method mentioned in Section 3.3.3. The result $a = 0.7002$ shows strong agreement that $\mathcal{H}$ is validated across all models based on converting the same outliers. This suggests that despite inconsistencies in how individual models integrate outliers over time, their validation of $\mathcal{H}$ remains broadly aligned. Accordingly, the average model $x = 0.80$ is a good consensus model regarding the validation of $\mathcal{H}$. For this reason, we proceed using the average model representation for our interpretability analysis.

In the next section, we examine whether writing style, beyond semantic similarity, helps explain why some outliers are eventually integrated into clusters, while others remain isolated.

## 4.4 Lexicon and Writing Style Analysis

As an attempt to explain the conversion of outliers into topics, we controlled for topical alignment among outliers to assess their influence on topic formation (see Subsection 3.3.4). For each word appearing in outlier documents, we computed the difference in average TF-IDF scores between those validating $\mathcal{H}$ and those not validating $\mathcal{H}$. Specifically, we used the lexical salience metric $\Delta\text{TFIDF}(w)$, as defined in (1), and its inverse. Among the top 20 words in each class, the mean difference was 0.0088 for outliers validating $\mathcal{H}$ and $-0.0126$ for those not validating $\mathcal{H}$. Both differences were statistically significant at the 0.05 level using the Kruskal–Wallis test.

A closer examination of the top 20 terms most prevalent among outliers validating $\mathcal{H}$ reveals words associated with institutional support for the project (e.g., "*cabinet*", "*total*", "*lobbying*", "*saclay*") or individuals endorsing it (e.g., "*brunelle*", "*nathalie*"). In contrast, the top 20 terms more prevalent among outliers *not* validating $\mathcal{H}$ include words reflecting opposition to the project (e.g., "*recours*", "*victoire*"), as well as references to activist groups (e.g., "*greenpeace*", "*militant*") and opposing figures (e.g., "*julliard*", "*jean*"). In both sets, the majority of these words were sta-

tistically distinctive.[2] These results suggest that conversion of outliers into topics is partly influenced by their alignment with dominant themes in the TP dataset, which is consistent with the role of semantic similarity in reinforcing or generating topical structure.

To evaluate our qualitative observation that the lexicon of outliers not validating $\mathcal{H}$ tends to be more subjectively framed or opinion-laden, we carried out a quantitative analysis to test this hypothesis. Specifically, we assessed whether lexical salience defined in (1) correlated with the degree of subjectivity or neutrality in the documents where each word occurred. For each word $w$, we computed the average subjectivity and neutrality scores across all documents $D_w$ in which it appeared:

$$\text{Subjectivity}(w) = \frac{1}{|D_w|} \sum_{d \in D_w} s(d) \quad (2)$$

$$\text{Neutrality}(w) = \frac{1}{|D_w|} \sum_{d \in D_w} n(d) \quad (3)$$

where $s(d)$ and $n(d)$ denote the subjectivity and neutrality scores of document $d$, computed using TextBlob (Loria et al., 2018) and VADER (Hutto and Gilbert, 2014), respectively.

We then computed Spearman's correlation between $\Delta\text{TFIDF}(w)$ and the subjectivity and neutrality scores of the corresponding documents. The analysis revealed a *moderate negative correlation* with subjectivity ($r = -0.223$, $p < 0.01$) and a *weak positive correlation* with neutrality ($r = 0.105$, $p < 0.01$). These patterns indicate that words more prominent among converted outliers tend to appear in more neutral, less subjective contexts and thus that outliers more likely to become topics are characterized by a lexicon that is more factual in nature.

To evaluate broader stylistic effects, we applied the stylometric framework of Terreau et al. (2021) to measure differences across eight core stylistic features between converted and non-converted outliers: *function words*, *punctuation marks*, *numbers*, *named entities*, *part-of-speech tags*, *structural features*, *lexical complexity indexes*, and *readability metrics*. Figure 3 summarizes the results for both main categories (Fig. 3a) and significant subfeatures (Fig. 3b). We omit a detailed analysis for

---

[2]Three words among outliers validating $\mathcal{H}$ —"*public*", "*direction*", and "*palaiseau*"—were not statistically significant, while only one word ("*ecole*") lacked significance among outliers not validating $\mathcal{H}$.

(a) Differences in mean frequencies for the eight main features.



(b) Differences in mean frequencies for subfeatures, based on observed significance in Figure 3a.

Figure 3: Differences for TP in the eight stylistic features and subfeatures from Terreau et al. (2021), between outliers validating $\mathcal{H}$ and outliers not validating $\mathcal{H}$. Statistical significance is measured using the Kruskal–Wallis test, with $*$ and $**$ indicating $p$ values $< 0.05$ and $< 0.01$, respectively.

significant differences in Numbers, as this feature consists solely of single-digit values (ranging from 0 to 9), making any further breakdown not directly interpretable.

At the level of the eight main features (see Figure 3a), statistically significant differences were observed only for Named Entities (NER), Structural features, and Numbers. Structural features and NER were less frequent in outliers validating $\mathcal{H}$ than in those not validating $\mathcal{H}$, whereas Numbers were more frequent in outliers validating $\mathcal{H}$. No significant differences were found for TAG, Punctuation, Letters, Indexes, or Function Words.

A closer examination of the significant subfeatures (Fig. 3b) shows that, for NER, names of persons and organizations appear significantly less often in outliers validating $\mathcal{H}$. No difference was

found for location markers. For Structural subfeatures, outliers validating $\mathcal{H}$ exhibit shorter sentences and words, fewer syllables per word, and higher average word frequency. No other structural subfeatures showed significant variation.

These stylistic differences observed for the average model may be explained by the fact that more structural features introduce complexity, and thus stylistic simplification may support the integration of outliers into topic clusters. Specifically, shorter and simpler text, with fewer named entities, may make it easier for the average model to associate such outliers with broader topic structures, thus facilitating the validation of $\mathcal{H}$. Conversely, a higher frequency of Numbers, particularly single-digit ones, may reflect more patterned or categorical language that also facilitates topic clustering. No clear effects were found for TAG, Punctuation, Letters, Indexes, or the remaining structural subfeatures.

## 5 Replication Study

### 5.1 English Dataset

To validate and generalize our findings, we used an existing larger English dataset of climate change news articles, *climate-news-db*.[3] This dataset originally comprised 27,877 news articles from global media outlets, spanning January 2015 to November 2024. To ensure topical consistency, we curated a focused subset of 312 articles, referred to as GHG, by filtering for content explicitly addressing Greenhouse Gas Emissions (GHG). Articles were selected based on the presence of the terms *"Greenhouse Gas"* or *"Greenhouse Emissions"*, and sampled across 20 monthly time windows between January 2022 and August 2023. For consistency, we retained only articles from major U.S.-based outlets (e.g., *The Washington Post*, *The New York Times*, *Fox News*, and *CNN*).

### 5.2 Topic-Based Clustering

We applied topic-based clustering to the body text of the GHG articles using 10D UMAP projections. With the exception of e5-base-v2, Table 2 shows that all nine models achieved strong silhouette scores, with both mean and median values at or above 0.5 (on a scale from −1 to 1). These results are slightly lower than, but broadly consistent with, those obtained for the TP dataset under the same configuration.

---

[3] https://www.climate-news-db.com

Figure 4: 2D Scatter plot of the `UMAP` 10D cumulative clustering obtained on `GHG` over nine time windows, using `e5-base-v2`. Outliers are indicated with black ×, topics in blue, green and yellow.

Table 2: `UMAP` 10D silhouette scores obtained on the `GHG` dataset for the body text of articles, sorted from best to worst.

| Model | Mean Silhouette Score |
|---|---|
| `e5-base-v2` | **0.5661** |
| `multilingual-e5-large` | 0.5490 |
| `all-MiniLM-L12-v2` | 0.5416 |
| `...-multi..-mpnet-base-v2` | 0.5387 |
| `xlm-roberta-large` | 0.5376 |
| `Solon-embeddings-large-0.1` | 0.5159 |
| `sentence-camembert-base` | 0.5092 |
| `all-roberta-large-v1` | 0.5044 |
| `distilbert-base-uncased` | 0.4998 |
| Mean | 0.5291 |
| Median | 0.5376 |

## 5.3 Outlier Behavior

Figure 5 shows the mean validation score per model for $\mathcal{H}$ on `GHG`. The results indicate a high average validation across models, with a mean score of 0.81. As expected, English-specialized models: `distilbert-base-uncased`, `e5-base-v2`, and `all-MiniLM-L12-v2`, achieve perfect validation (1.0), followed by `all-roberta-large-v1` (0.85). Among French-specialized models, `sentence-camembert-base` performs more weakly (0.58), as anticipated, while the perfect score of `Solon-embeddings-large-0.1` (1.0) is less expected. Multilingual models show mixed results: `paraphrase-multilingual-mpnet-base-v2` and `xlm-roberta-large` perform poorly (both 0.41), while `multilingual-e5-large` again achieves perfect validation. The distribution of scores appears bimodal: five models achieve perfect validation,

while the remaining four show moderate to low scores. This sharp divide may reflect potential overfitting among English-specialized models that integrate all outliers into topics.



Figure 5: Mean number of outliers per model that validate prediction $\mathcal{H}$ on `GHG` by converting into topic inlier at some time point (specific to each model). Each colored bar represents the mean for each model.

This consistency in temporal dynamics (see Appendix A.4.1 for a detailed time-window analysis) aligns with the high average validation score of 0.81 (Figure 5). Most models follow a similar pattern: strong early outlier-to-topic conversion, reduced integration in mid-phases, and stabilization with persistent outliers. While some models, particularly multilingual ones and `sentence-camembert-base`, show greater fluctuation, the overall trend supports $\mathcal{H}$. As in the Pilot experiment, we computed inter-agreement across models with respect to $\mathcal{H}$, using the rescaling method of Icard et al. (2024). Again, the result $a = 0.6783$ strongly supports that models validate $\mathcal{H}$ based on converting the exact same outliers. The average model $x = 0.81$ is then a good consensus model regarding the validation of $\mathcal{H}$.

## 5.4 Lexicon and Writing Style Analysis

As part of our interpretability analysis, we sought to understand why some outliers aligned with topics while others did not. We first examined the top 20 words with the highest $\Delta$TFIDF scores in outliers validating $\mathcal{H}$ compared to those not validating it, and vice versa. As defined in (1), $\Delta$TFIDF$(w)$ captures the difference in average TFIDF scores for word $w$ between the two outlier classes. The mean difference was 0.0031 for (1), and 0.0023 for the reverse. Neither difference was statistically significant (Kruskal–Wallis test, $p > 0.05$), suggesting that thematic lexical content does not meaningfully distinguish the two outlier classes in `GHG`. How-

Figure 6: Differences for GHG in the eight stylistic features from Terreau et al. (2021), between outliers validating $\mathcal{H}$ and outliers not validating $\mathcal{H}$. Statistical significance is measured using the Kruskal–Wallis test, with * and ** indicating $p$ values $< 0.05$ and $< 0.01$, respectively.

ever, this finding does not rule out the possibility that stylistic or other non-topical lexical and linguistic features influence outlier conversion.

To address this gap, we analyzed the differences in stylistic characteristics between the two outlier classes, using the framework proposed by Terreau et al. (2021). The results for GHG are given in Figure 6 for the eight main features. We do not provide a detailed analysis of Function words and Letters, as Letters consist solely of single-character values (ranging from A to Z), and Function Words gain significance from their overall distribution rather than their individual occurrences, making a further breakdown not directly interpretable.

Among the eight features, significant frequency differences were found only for Function words and Letters, which were notably less frequent in outliers verifying $\mathcal{H}$ compared to outliers not verifying $\mathcal{H}$. This may be explained by the fact that function words (e.g., prepositions, conjunctions) and letters (e.g., A, B, C) lack semantic content, so their reduction helps the average model recognize topics in outliers and validate $\mathcal{H}$. In contrast, Indexes, Numbers, NER, Punctuation, TAG, and Structural features do not appear to have a particular effect on this recognition.

## 6  Discussion

We observed consistent outlier-to-topic conversion across two linguistically distinct datasets, confirming that the phenomenon generalizes. Validation of $\mathcal{H}$ is robust across topic domains (social respon-

sibility and climate change), languages (French and English), and dataset sizes (102 and 312 articles), with a stable mean score around 0.80. Inter-model agreement remains high (with $a = 0.7002$ for French, $a = 0.6783$ for English), suggesting that topic-based clustering reliably integrates outliers under varied conditions.

In lexical analysis, TF-IDF differences between converted and non-converted outliers were significant in TP but not in GHG. In TP, converted outliers were more strongly associated with lower subjectivity and higher lexical neutrality. This reflects a structural difference: TP focuses on a defined controversy with a polarized lexicon, while GHG likely follows a more neutral, report-oriented style, as it was not curated under controversy criteria.

The stylistic features analysis revealed that writing style has a significant impact on the conversion of outliers into topics, though the relevant features differ by language. In TP, conversion is influenced by structural features, named entities, and numbers; in GHG, by function words and letter distributions. This suggests that embedding models rely on language-specific stylistic cues when integrating outliers.

These differences align with model training: French-trained models perform better on TP, English-trained ones on GHG, while multilingual models show mixed results, reflecting their training data (see Table 3 in Appendix A.4.1 for details).

## 7  Conclusion

Our findings demonstrate that outlier-to-inlier conversion is a consistent mechanism in topic emergence within cumulative, density-based clustering frameworks. The effect is robust across nine language models, two typologically distinct languages, and datasets with varying topical scope. In the French dataset (TP), focused on a well-defined controversy, average model validation reached 0.80; in the English dataset (GHG), covering broader climate discourse, the score was similarly high at 0.81. Inter-model agreement exceeded 0.65 in both cases, indicating stable clustering dynamics across architectures and domains.

Future work will distinguish between outliers that act as precursors to new topics and those that reinforce existing structures. We aim to quantify their predictive value and examine their temporal behavior across phases of topic development.

We also plan to scale our analysis to larger and

more heterogeneous corpora, particularly in domains where informational risks, such as discursive conflict and disinformation, are likely to emerge or escalate. In parallel, we will evaluate alternative clustering algorithms with integrated outlier detection (e.g., OPTICS) and broaden our assessment across additional model architectures. These extensions aim to test the generality and deepen the explanatory power of our findings.

## Limitations

This study was designed as a controlled pilot to explore the predictive role of outliers in topic emergence under well-defined experimental conditions. Although the number of raw articles was relatively limited (102 in French and 312 in English), each document was processed with nine distinct language models, resulting in 918 French and 2,808 English data points. This mitigated the limitations typically associated with small corpus sizes.

High inter-model agreement ($a = 0.7002$ for French and $a = 0.6783$ for English) and consistent clustering quality (silhouette scores of 0.61 and 0.52, respectively) further support that the results are robust within the bounds of this setup.

The decision to prioritize depth over breadth at the expense of dataset size was deliberate: it enabled the construction of a high-quality, manually curated corpus with full-text availability, temporal continuity (i.e., no temporal gaps), and source diversity. This design helped control for confounding factors such as incomplete timelines and uneven topic coverage, which often affect large-scale datasets whose compilation processes are not fully transparent.

While these constraints were necessary to ensure experimental clarity and interpretability, they naturally limit the generalizability of the findings. Future work will scale the analysis to larger corpus of news articles to test its applicability in more complex and dynamic information environments.

## Ethics Statement

Our research adheres to the ethical principles of open science, transparency, and sustainability. We ensure reproducibility by making our code accessible in a dedicated GitHub repository, with data and results available upon request. We comply with intellectual property and data protection regulations by sharing only vector embeddings generated by language models. This approach aligns with the principles of 'transformative fair use". We promote AI transparency by contributing to the interpretability of language models, supporting the responsible and explainable use of these models. To support efficiency and sustainability, we prioritize the use of small, open-source language models.

## Declaration of contribution

EZ, BI, and JGG conceptualized the research problem and designed the experiments. EZ managed the data collection process. AB and LS managed data cleaning and annotation. EZ was responsible for the technical aspects: coding, model selection, and building the experimental framework. EZ, BI, GB, and JGG analyzed and discussed the results. EZ and BI wrote the paper, which all authors read and revised together. EZ and BI share first authorship. Correspondence: evangelia.zve@lip6.fr, benjamin.icard@lip6.fr, jean-gabriel.ganascia@lip6.fr.

## References

Fuad Alattar and Khaled Shaalan. 2021. Emerging research topic detection using filtered-lda. *AI*, 2(4):578–599.

Mihael Ankerst, Markus M Breunig, Hans-Peter Kriegel, and Jörg Sander. 1999. Optics: Ordering points to identify the clustering structure. *ACM Sigmod record*, 28(2):49–60.

Daniel Atzberger, Tim Cech, Willy Scheibel, Matthias Trapp, R. Richter, J. Dollner, and Tobias Schreck. 2023. Large-scale evaluation of topic models and dimensionality reduction methods for 2d text spatialization. *DOI: 10.1109/TVCG.2023.3326569*. Accessed: January 2025.

Olusola Babalola, Bolanle Ojokoh, and Olutayo Boyinbode. 2024. Comprehensive evaluation of lda, nmf, and bertopic's performance on news headline topic modeling. *Journal of Computing Theories and Applications*, 2(2):268–289.

David M Blei, Andrew Y Ng, and Michael I Jordan. 2003. Latent dirichlet allocation. *Journal of machine Learning research*, 3(Jan):993–1022.

Allaa Boutaleb, Jerome Picault, and Guillaume Grosjean. 2024. Bertrend: Neural topic modeling for emerging trends detection. *arXiv preprint arXiv:2411.05930*.

Markus M Breunig, Hans-Peter Kriegel, Raymond T Ng, and Jörg Sander. 2000. Lof: Identifying density-based local outliers. In *Proceedings of the 2000 ACM SIGMOD international conference on Management of data*, pages 93–104. ACM.

Ricardo JGB Campello, Davoud Moulavi, Arthur Zimek, and Jörg Sander. 2015. Hierarchical density estimates for data clustering, visualization, and outlier detection. *ACM Transactions on Knowledge Discovery from Data (TKDD)*, 10(1):1–51.

Yubo Chen, Liheng Xu, Kang Liu, and Jun Zhao. 2016. Event detection with burst information networks. In *Proceedings of the 26th International Conference on Computational Linguistics (COLING)*, pages 3317–3327.

Maarten Grootendorst. 2022. Bertopic: Neural topic modeling with pre-trained language models. *Journal of Machine Learning Applications*, 17(3):45–62.

John A Hartigan and Manchek A Wong. 1979. Algorithm as 136: A k-means clustering algorithm. *Journal of the Royal Statistical Society: Series C (Applied Statistics)*, 28(1):100–108.

Amos Hoyle, Pratyusha Goel, Cynthia Phillips, Jordan Boyd-Graber, and Philip Resnik. 2022. Is automated topic model evaluation broken? the incoherence of coherence. In *Findings of the Association for Computational Linguistics: ACL 2022*, pages 478–493.

Clayton Hutto and Eric Gilbert. 2014. Vader: A parsimonious rule-based model for sentiment analysis of social media text. In *Proceedings of the international AAAI conference on web and social media*, volume 8, pages 216–225.

Benjamin Icard, François Maine, Morgane Casanova, Géraud Faye, Julien Chanson, Guillaume Gadek, Ghislain Atemezing, François Bancilhon, and Paul Égré. 2024. A multi-label dataset of french fake news: Human and machine insights. In *Proceedings of the 2024 Joint International Conference on Computational Linguistics, Language Resources and Evaluation (LREC-COLING 2024)*, pages 812–818.

Sukhwan Jung, Rituparna Datta, and Aviv Segev. 2020. Identification and prediction of emerging topics in news media. In *IEEE International Conference on Big Data*.

Charu Karakkaparambil, Mayank Nagda, Nooshin Haji Ghassemi, Marius Kloft, and Sophie Fellenz. 2024. Evaluating dynamic topic models. In *Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (ACL)*.

J Peter Kincaid, Robert P Fishburne Jr, Richard L Rogers, and Brad S Chissom. 1975. Derivation of new readability formulas (automated readability index, fog count and flesch reading ease formula) for navy enlisted personnel.

Mario Köppen. 2000. The curse of dimensionality. In *5th online world conference on soft computing in industrial applications (WSC5)*, volume 1, pages 4–8.

Caitlin Doogan Poet Laureate, Wray Buntine, and Henry Linger. 2023. A systematic review of the use of topic models for short text social media analysis. *Artificial Intelligence Review*, 56(12):14223–14255.

Haein Lee, Seon Hong Lee, Kyeo Re Lee, and Jang Hyun Kim. 2023. Esg discourse analysis through bertopic: comparing news articles and academic papers. *Computers, Materials & Continua*, 75(3):6023–6037.

Fei Tony Liu, Kai Ming Ting, and Zhi-Hua Zhou. 2008. Isolation forest. In *2008 eighth ieee international conference on data mining*, pages 413–422. IEEE.

Steven Loria et al. 2018. textblob documentation. *Release 0.15*, 2(8):269.

Leland McInnes, John Healy, Steve Astels, et al. 2017. hdbscan: Hierarchical density based clustering. *J. Open Source Softw.*, 2(11):205.

Leland McInnes, John Healy, and James Melville. 2018. Umap: Uniform manifold approximation and projection for dimension reduction. *arXiv preprint arXiv:1802.03426*.

Niklas Muennighoff, Nouamane Tazi, Loïc Magne, and Nils Reimers. 2022. Mteb: Massive text embedding benchmark. *arXiv preprint arXiv:2210.07316*.

Shahzad Qaiser and Ramsha Ali. 2018. Text mining: use of tf-idf to examine the relevance of words to documents. *International journal of computer applications*, 181(1):25–29.

Ketan Rajshekhar Shahapure and Charles Nicholas. 2020. Cluster quality analysis using silhouette score. In *2020 IEEE 7th international conference on data science and advanced analytics (DSAA)*, pages 747–748. IEEE.

Claude Elwood Shannon. 1948. A mathematical theory of communication. *The Bell system technical journal*, 27(3):379–423.

Enzo Terreau, Antoine Gourru, and Julien Velcin. 2021. Writing style author embedding evaluation. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 84–93. Association for Computational Linguistics.

H. Wold. 1987. Principal component analysis. *Technometrics*, 38(3):235–238.

Lili Yao, Yue Zhang, and Xisen Wang. 2021. Dynamic term frequency analysis for topic shift detection. *ACL*.

Tuukka Ylä-Anttila, Veikko Eranti, and Anna Kukkonen. 2022. Topic modeling for frame analysis: A study of media debates on climate change in india and usa. *Global Media and Communication*, 18(1):91–112.

C Udny Yule. 2014. *The statistical study of literary vocabulary*. Cambridge University Press.

394

## A Appendix

### A.1 Supplementary Materials

The code and visualizations supporting this paper are available at: https://github.com/evangeliazve/outliers-to-topics-icnlsp. The datasets and experimental results can be provided upon request. The repository includes Python scripts for reproducing our experiments, as well as statistical analyses and visualizations corresponding to key figures and tables in the paper. The BERTopic framework is documented at: https://maartengr.github.io/BERTopic/. Further details on HDBSCAN can be found in its official documentation: https://hdbscan.readthedocs.io/en/latest/, and information on UMAP dimensionality reduction is available at: https://umap-learn.readthedocs.io/en/latest/basic_usage.html. For TF-IDF, we used the TfidfVectorizer from scikit-learn: https://scikit-learn.org/stable/modules/generated/sklearn.feature_extraction.text.TfidfVectorizer.html. The sentiment analysis tools employed in this study are TextBlob (https://textblob.readthedocs.io/en/dev/index.html) and VADER (https://github.com/cjhutto/vaderSentiment)..

### A.2 Models

Table 3 presents the nine sentence embedding models used in our experiments for topic-based clustering, detailing their underlying architectures, embedding dimensionality, language coverage, and model sizes.

### A.3 Pilot Study Appendix

#### A.3.1 Detailed Silhouette Scores

In this appendix, we provide detailed results from the pilot study evaluating the effectiveness of different sentence embedding models for topic-based clustering on the TP dataset. Specifically, Table 4 reports the mean silhouette scores obtained for each model under varying dimensionality reductions (2D, 3D, 5D, and 10D using UMAP) and different text sample types (headline, body, and combined text). These results offer insights into how model selection, dimensionality, and text granularity impact clustering quality.

#### A.3.2 Validation or invalidation of $\mathcal{H}$ per model over different time windows for TP

For a more detailed examination of model variations, both across and within models, Table 5 presents the validation or invalidation of $\mathcal{H}$ per model over different periods of cumulative clustering. Among French models, Solon-embeddings-large-0.1 is fully consistent, achieving complete integration early, while sentence-camembert-base shows non-monotony, with conversion dropping from 95.83% to 50% and some persistent outliers. English models exhibit pronounced inconsistency: e5-base-v2 weakens over time (68.75% to 36.67%), all-MiniLM-L12-v2 and all-roberta-large-v1 show fluctuating progress despite strong early conversion (66.67% and 85.42%), and distilbert-base-uncased remains fully stable with no outliers through the whole. Multilingual models vary widely, with paraphrase-multilingual-mpnet-base-v2 and xlm-roberta-large starting strong (85.42%, 64.58%) but leaving substantial outliers later (12.74%, 42.15%), while multilingual-e5-large follows an unstable trajectory, declining from 83.33% to 48.15%.

Across models, a general pattern emerges: strong early conversion of outliers into topic inliers, slowing integration in the mid-phase, and eventual stabilization with persistent outliers in 2023. Early clustering is largely consistent, with conversion rates ranging from 64.58% (xlm-roberta-large) to 100% (Solon-embeddings-large-0.1) in 2020-11. By the mid-phase (2021-07), some models, like all-MiniLM-L12-v2 (73.91%) and all-roberta-large-v1 (71.43%), sustain moderate integration, while others, like e5-base-v2 (42.86%), decline. Late-stage variations are more pronounced, with paraphrase-multilingual-mpnet-base-v2 retaining 46.47% of outliers as topic inliers, while xlm-roberta-large and e5-base-v2 drop to 26.67% and 36.67%, respectively. sentence-camembert-base, despite an early peak (95.83%), declines to 50.00%.

#### A.3.3 Top 10 Distinguishing Terms Based on TF-IDF Differences Between Outliers Validating and Not Validating $\mathcal{H}$

Table 6 lists the top 10 terms whose TF-IDF scores most strongly differentiate outliers that validate hypothesis $\mathcal{H}$ from those that do not, highlighting

| Model | Architecture | Dimensions | Language | Parameters |
|---|---|---|---|---|
| Solon-embeddings-large-0.1 | RoBERTa | 1024 | French | 560M |
| sentence-camembert-base | CamemBERT | 768 | | 111M |
| all-roberta-large-v1 | RoBERTa | 1024 | English | 355M |
| e5-base-v2 | E5 | 768 | | 109M |
| distilbert-base-uncased | DistilBERT | 768 | | 67M |
| all-MiniLM-L12-v2 | MiniLM | 384 | | 33.4M |
| xlm-roberta-large | XLM-RoBERTa | 1024 | Multilingual | 561M |
| multilingual-e5-large | E5 | 1024 | | 560M |
| paraphrase-multilingual-mpnet-base-v2 | MPNet | 768 | | 278M |

Table 3: Description of the nine sentence embedding models used to conduct the topic-based clustering experiments.

| Model | UMAP 2D | | | UMAP 3D | | | UMAP 5D | | | UMAP 10D | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | Headline | Body | All | Headline | Body | All | Headline | Body | All | Headline | Body | All |
| multilingual-e5-large | 0.6235 | 0.6002 | 0.5914 | 0.6121 | 0.5480 | 0.5713 | 0.6020 | 0.5481 | 0.5692 | 0.6065 | 0.5519 | 0.5689 |
| e5-base-v2 | 0.6133 | 0.5556 | 0.4668 | 0.5718 | 0.5627 | 0.4671 | 0.5580 | 0.5479 | 0.5030 | 0.5592 | 0.5350 | 0.4846 |
| sentence-camembert-base | 0.6120 | 0.5616 | 0.5994 | 0.6083 | 0.5791 | 0.6302 | 0.5934 | 0.5877 | 0.6354 | 0.5990 | 0.5850 | 0.6167 |
| all-MiniLM-L12-v2 | 0.6039 | 0.5858 | 0.5465 | 0.5570 | 0.6197 | 0.5243 | 0.5702 | 0.4962 | 0.5608 | 0.5654 | 0.5846 | 0.5349 |
| Solon-embeddings-large-0.1 | 0.5573 | 0.6340 | 0.6497 | 0.6056 | 0.6351 | 0.6031 | 0.5660 | 0.6153 | 0.5778 | 0.5772 | 0.6694 | 0.5553 |
| xlm-roberta-large | 0.5416 | 0.3729 | 0.3294 | 0.4996 | 0.3812 | 0.3226 | 0.5348 | 0.3694 | 0.3848 | 0.4941 | 0.4802 | 0.4424 |
| all-roberta-large-v1 | 0.5294 | 0.6701 | 0.5862 | 0.5427 | 0.6255 | 0.6121 | 0.5536 | 0.6361 | 0.6040 | 0.5525 | 0.6258 | 0.5759 |
| ..-multilingual-mpnet-base-v2 | 0.5259 | 0.6062 | 0.7324 | 0.5221 | 0.5918 | 0.6429 | 0.5517 | 0.5977 | 0.6754 | 0.5391 | 0.5923 | 0.6865 |
| distilbert-base-uncased | 0.4872 | 0.7907 | 0.8535 | 0.5232 | 0.9233 | 0.8509 | 0.4413 | 0.9575 | 0.8670 | 0.3670 | 0.9373 | 0.8895 |
| Mean | 0.5660 | 0.5975 | 0.5945 | 0.5603 | 0.6074 | 0.5816 | 0.5523 | 0.5951 | 0.6008 | 0.5400 | 0.6180 | 0.5993 |
| Median | 0.5588 | 0.6056 | 0.5929 | 0.5692 | 0.5984 | 0.5718 | 0.5544 | 0.6029 | 0.6040 | 0.5417 | 0.6183 | 0.5756 |

Table 4: Mean silhouette scores per model, dimensionality and text samples types obtained on dataset TP.

key lexical features associated with each group.

## A.4 Replication Study Appendix

### A.4.1 Validation or invalidation of $\mathcal{H}$ per model over different periods for GHG

For a detailed examination of model variations, Table 7 presents the validation or invalidation of $\mathcal{H}$ for each model over different periods of cumulative clustering. Among English models, distilbert-base-uncased, e5-base-v2, and all-MiniLM-L12-v2 show complete consistency, achieving full integration early. all-roberta-large-v1 follows a steady trajectory, with conversion decreasing slightly from 93.62% to 88.17%. Among French models, sentence-camembert-base, a French model, shows instability, with a conversion fluctuating from 46.43% to 58.14% before dropping to 43.18%. Solon-embeddings-large-0.1, despite being a French model, integrates all outliers early, aligning with its high absolute validation score. Multilingual models exhibit mixed behaviors., with multilingual-e5-large achieving full integration like English models, while paraphrase-multilingual-mpnet-base-v2 and xlm-roberta-large retain substantial outliers (with 40.70% and 43.91%, respectively). multilingual-mpnet-base-v2 initially increases its conversion (26.32% to 38.10%) before stabilizing. xlm-roberta-large exhibits a downward trend, with

conversion dropping from 45.00% to 22.00%.

That said, trends across models reveal a broadly consistent trajectory: high early conversion of outliers into topic inliers (ranging from 26.32% to 100% in 2022-10), followed by a mid-phase slowdown with moderate-to-low integration (10.95%–65.71% in 2023-02), and eventual stabilization with persistent outliers in the final stage (10.25%–43.91%). Most models adhere to this pattern, with strong early conversion seen in all-roberta-large-v1 (93.62%) and e5-base-v2 (100%), followed by a gradual decline in mid-phase integration for models like sentence-camembert-base (fluctuating from 46.43% to 58.14%) and multilingual-mpnet-base-v2 (increasing from 26.32% to 38.10%). By the final stage, outlier retention converges to similar rates across models, such as sentence-camembert-base stabilizing at 25.64% and xlm-roberta-large retaining 43.91% of outliers.

### A.4.2 Top 10 Distinguishing Terms Based on TF-IDF Differences Between Outliers Validating and Not Validating $\mathcal{H}$

Table 8 lists the top 10 terms whose TF-IDF scores most strongly differentiate outliers that validate hypothesis $\mathcal{H}$ from those that do not, highlighting key lexical features associated with each group.

| Model | Measures | Time | | | |
|---|---|---|---|---|---|
| | | **2020-11 (50%)** | **2021-07 (70%)** | **2023-09 (90%)** | **Remaining (100%)** |
| `Solon-..-large-0.1` | Nb Outliers / All Articles at $t$ | 48/48 | 8/69 | 0/100 | 0/102 |
| | % Becoming Inliers at $(t+n)$ | 100% | 100% | 0.00 | Converted on 2022-01 |
| `..-multi..-mpnet-..` | Nb Outliers / All Articles at $t$ | 48/48 | 69/69 | 15/100 | 13/102 |
| | % Becoming Inliers at $(t+n)$ | 85.42% | 84.06% | 46.47% | - |
| `sentence-camembert-..` | Nb Outliers / All Articles at $t$ | 48/48 | 33/69 | 8/100 | 4/102 |
| | % Becoming Inliers at $(t+n)$ | 95.83% | 90.91% | 50.00% | - |
| `multi..-e5-large` | Nb Outliers / All Articles at $t$ | 48/48 | 25/69 | 27/100 | 25/102 |
| | % Becoming Inliers at $(t+n)$ | 83.33% | 64.00% | 48.15% | - |
| `xlm-roberta-large` | Nb Outliers / All Articles at $t$ | 48/48 | 39/69 | 30/100 | 43/102 |
| | % Becoming Inliers at $(t+n)$ | 64.58% | 64.10% | 26.67% | - |
| `all-MiniLM-L12-v2` | Nb Outliers / All Articles at $t$ | 48/48 | 69/69 | 16/100 | 26/102 |
| | % Becoming Inliers at $(t+n)$ | 66.67% | 73.91% | 12.50% | - |
| `all-roberta-large-v1` | Nb Outliers / All Articles at $t$ | 48/48 | 21/69 | 10/100 | 12/102 |
| | % Becoming Inliers at $(t+n)$ | 85.42% | 71.43% | 30.00% | - |
| `distil..-base-uncased` | Nb Outliers / All Articles at $t$ | 0/48 | 0/69 | 0/100 | 0/102 |
| | % Becoming Inliers at $(t+n)$ | 0.00% | 0.00% | 0.00% | Converted on 2020-06 |
| `e5-base-v2` | Nb Outliers / All Articles at $t$ | 48/48 | 21/69 | 30/100 | 41/102 |
| | % Becoming Inliers at $(t+n)$ | 68.75% | 42.86% | 36.67% | - |

Table 5: Proportion of outliers converting to clusters in TP, for each model and along four time windows.

| Word | $\Delta$TFIDF$(w)$ | $\Delta$Occ$(w)$ | Word | $\Delta$TFIDF$(w)$ | $\Delta$Occ$(w)$ |
|---|---|---|---|---|---|
| cabinet | 0.0122* | 93 | totalenergies | -0.0328** | -28 |
| total | 0.0119* | 2613 | recours | -0.0185** | -14 |
| brunelle | 0.0106* | 136 | greenpeace | -0.0173** | -265 |
| nathalie | 0.0104* | 139 | victoire | -0.0155** | -6 |
| lobbying | 0.0103** | 122 | ecole | -0.0143 | -162 |
| public | 0.0098 | 428 | julliard | -0.0129** | -14 |
| direction | 0.0097 | 563 | jean | -0.0126** | -20 |
| palaiseau | 0.0095 | 60 | décision | -0.0124** | -127 |
| saclay | 0.0089* | 740 | conseil | -0.0116** | -626 |
| quartier | 0.0086* | 40 | militant | -0.0112** | -47 |

Table 6: Top 10 absolute values of $\Delta$TFIDF$(w)$ for TP. Words with positive values are more characteristic of converted outliers ($\mathcal{H}$), and those with negative values are more typical of non-converted outliers (not $\mathcal{H}$). Statistical significance is based on the Kruskal-Wallis test; * and ** indicate $p$-values $< 0.05$ and $< 0.01$, respectively. $\Delta$Occ$(w)$ indicates the difference in word occurrence counts between the two groups.

| Model | Measures | Time | | | |
|---|---|---|---|---|---|
| | | **2022-10 (50%)** | **2023-02 (70%)** | **2023-06 (90%)** | **Remaining (100%)** |
| Solon-embeddings-large-0.1 | Nb Outliers / All Articles at $t$ | 79/79 | 18/105 | 96/236 | 0/312 |
| | % Becoming Inliers at $(t+n)$ | 100% | 100% | 100% | Converted on 2023-07 |
| ...-multi..-mpnet-base-v2 | Nb Outliers / All Articles at $t$ | 19/79 | 21/105 | 81/236 | 127/312 |
| | % Becoming Inliers at $(t+n)$ | 26.32% | 38.10% | 33.33% | - |
| sentence-camembert-base | Nb Outliers / All Articles at $t$ | 28/79 | 43/105 | 88/236 | 80/312 |
| | % Becoming Inliers at $(t+n)$ | 46.43% | 58.14% | 43.18% | - |
| multi..-e5-large | Nb Outliers / All Articles at $t$ | 23/79 | 26/105 | 49/236 | 0/312 |
| | % Becoming Inliers at $(t+n)$ | 100% | 100% | 100% | Converted on 2023-07 |
| xlm-roberta-large | Nb Outliers / All Articles at $t$ | 20/79 | 60/105 | 76/236 | 137/312 |
| | % Becoming Inliers at $(t+n)$ | 45.00% | 45.00% | 22.00% | - |
| all-MiniLM-L12-v2 | Nb Outliers / All Articles at $t$ | 79/79 | 69/105 | 90/236 | 0/312 |
| | % Becoming Inliers at $(t+n)$ | 100% | 100% | 100% | Converted on 2023-07 |
| all-roberta-large-v1 | Nb Outliers / All Articles at $t$ | 47/79 | 40/105 | 93/236 | 32/312 |
| | % Becoming Inliers at $(t+n)$ | 93.62% | 92.50% | 88.17% | - |
| distilbert-base-uncased | Nb Outliers / All Articles at $t$ | 42/79 | 33/105 | 87/236 | 0/312 |
| | % Becoming Inliers at $(t+n)$ | 100% | 100% | 100% | Converted on 2023-07 |
| e5-base-v2 | Nb Outliers / All Articles at $t$ | 13/79 | 27/105 | 58/236 | 0/312 |
| | % Becoming Inliers at $(t+n)$ | 100% | 100% | 100% | Converted on 2023-07 |

Table 7: Proportion of outliers converting to clusters in GHG, for each model and along four time windows.

| Word | $\Delta$TFIDF($w$) | $\Delta$Occ($w$) | Word | $\Delta$TFIDF($w$) | $\Delta$Occ($w$) |
|---|---|---|---|---|---|
| climate | 0.0067 | 17851 | amazon | -0.0034 | -98 |
| report | 0.0051* | 2656 | pakistan | -0.0033 | -130 |
| degree | 0.0035 | 2576 | china | -0.0031 | -480 |
| said | 0.0035 | 9134 | child | -0.0027 | -227 |
| bill | 0.0033* | 804 | thunberg | -0.0024 | -63 |
| company | 0.0032 | 2577 | reactor | -0.0023 | -65 |
| would | 0.0031 | 3668 | protest | -0.0023 | -87 |
| republican | 0.0030 | 728 | soil | -0.0023 | -185 |
| energy | 0.0030 | 5651 | granholm | -0.0023 | -51 |
| nice | 0.0029 | 895 | art | -0.0023 | -172 |

Table 8: Top 10 absolute values of $\Delta$TFIDF($w$) for GHG. Words with positive values are more characteristic of converted outliers ($\mathcal{H}$); words with negative values are more typical of non-converted outliers (not $\mathcal{H}$). Statistical significance is based on the Kruskal-Wallis test; * indicates $p < 0.05$. $\Delta$Occ($w$) shows the difference in word frequency between the two groups.

# From Context to Emotion: Leveraging LLMs for Recognizing Implicit Emotions

**Hanane Boutouta, Abdelaziz Lakhfif, Ferial Senator,** and **Chahrazed Mediani**
Artificial Intelligence laboratory (AI-lab), Department of Computer Science,
Setif 1 University, Ferhat ABBAS, Setif, Algeria
{hanane.boutouta, abdelaziz.lakhfif, ferial.senator, chahrazed.mediani}@univ-setif.dz

## Abstract

Implicit Emotion Recognition (IER) is a challenging task in Natural Language Processing (NLP), as it requires identifying emotions that are not directly expressed through explicit emotion words but must be inferred from contextual, situational, or linguistic cues. With the rapid progress of Large Language Models (LLMs), new opportunities have emerged for tackling such complex language understanding tasks. In this work, we investigate the effectiveness of two different architectures of LLMs for IER: masked language models, including BERT and RoBERTa, and causal language models, represented by ChatGPT. We fine-tuned BERT and RoBERTa on benchmark IER datasets, while we evaluated ChatGPT in a zero-shot setting to assess its ability to generalize without task-specific training. Our experiments on the ISEAR and IEST datasets show that fine-tuned masked language models perform strongly on the IER task. At the same time, ChatGPT achieves promising results in zero-shot scenarios, highlighting its potential for emotion recognition tasks with limited or no labeled data.

## 1 Introduction

Text-based Emotion Recognition (ER) is a fundamental research area in Natural Language Processing (NLP). In recent years, this field has seen important advancements due to increased human-computer interaction, as well as the rapid growth of online social media (Bisogni et al., 2023). ER can be classified into explicit ER (EER) and implicit ER (IER), depending on whether explicit emotional words emerge in the text (Kusal et al., 2021). Different from EER, where emotional words (e.g., happy, angry) occur in the text, in IER, emotions must be inferred from linguistic cues such as contextual descriptions, metaphorical expressions, or situational events without any explicit emotional expression (Klinger et al., 2018). Implicit emotions often

require deep semantic understanding to interpret subtle cues, such as sarcasm (Perfect, just what I needed) (Zhu et al., 2025), ambiguous statements (e.g., She looked out the window as the train pulled away, which could imply sadness, longing, or even relief) (Orizu, 2018), or behavioral context (e.g., They all left without me). This makes IER particularly challenging due to subjectivity, cultural variability, and strong dependence on context.

Researchers have proposed several ER approaches, including lexicon-based, machine learning, and deep learning methods. Most of these approaches primarily focus on extracting explicit emotions, whereas recognizing implicit emotions poses a greater challenge, as it demands sophisticated techniques capable of accurately interpreting context and deeply understanding nuanced linguistic patterns.

Recent advancements in Large Language Models (LLMs) have revolutionized NLP by achieving state-of-the-art performance across a wide range of tasks, such as question answering (Goar et al., 2023), machine translation (Hendy et al., 2023), and sentiment analysis (Ding et al., 2022). These models have also demonstrated remarkable capabilities in comprehending, interpreting, and recognizing human emotions (Banimelhem and Amayreh, 2023; Lee et al., 2024). LLMs, trained on large-scale and extensive corpora, have demonstrated a deep understanding of linguistic patterns, contextual dependencies, and even some aspects of world knowledge, enabling them to infer meaning and emotion from text based on the surrounding context in ways that were previously unattainable (Hong et al., 2024; Buscemi and Proverbio, 2024). These capabilities may be particularly useful for tasks like IER, where emotions are not explicitly stated but must be inferred from subtle linguistic cues, situational context, or background knowledge. Unlike traditional methods that rely heavily on explicit emotional keywords or rule-based systems, LLMs

leverage their transformer-based architecture, contextual embeddings, and pretrained knowledge to analyze the interplay of words and sentences, to decode emotional tones and comprehend the complexity of emotions.

In this study, we aim to explore the effectiveness of LLMs, specifically BERT, RoBERTa, and ChatGPT, in the task of IER. We assess the performance of fine-tuned, encoder-based models, including BERT and RoBERTa architectures, to evaluate their suitability and effectiveness for IER. Furthermore, we investigate ChatGPT's capabilities in a zero-shot learning setting to determine its ability to generalize to IER without task-specific fine-tuning. This approach highlights its potential for applications where labeled data is limited or unavailable.

This comparison provides insight into the strengths and weaknesses of these LLMs in capturing implicit emotional cues, contributing to a deeper understanding of their real-world applicability. The main contributions of this paper are summarized as follows:

- We conduct a comparative analysis of two distinct paradigms for IER:
    1. Fine-tuned masked language models (BERT and RoBERTa), and
    2. Zero-shot prompting using a causal language model, specifically ChatGPT.

- We fine-tune BERT and RoBERTa on labeled emotion datasets to evaluate their task-specific performance in recognizing implicit emotions.

- We evaluate the performance of ChatGPT to generalize to the IER task, in a zero-shot setting without the need for task-specific fine-tuning or additional training.

- We provide empirical evidence on the effectiveness, limitations, and generalization capabilities of ChatGPT in contrast to traditional fine-tuned models.

The remainder of this paper is structured as follows: Section 2 presents a concise review of related work on IER and recent developments in LLMs. Section 3 introduces the datasets used in our experiments. Section 4 provides some details about the experiments' setup, including the models employed and the different methodological approaches used. Section 5 presents and discusses the results of our experiments. Finally, Section 6 concludes the paper and outlines potential directions for future research.

## 2 Related Work

IER has emerged as a complex and less explored task within the field of NLP. Unlike EER, which relies on identifying overt emotion words, IER requires understanding contextual and semantic cues to infer emotional states. This task has been addressed using various approaches, including rule-based, classical machine learning, deep learning, and, more recently, transformer-based approaches (Alswaidan and Menai, 2020).

Early efforts relied on knowledge-based and rule-based approaches. For instance, EmotiNet linked events to emotions through commonsense knowledge, while cognitive-theory-inspired rules attempted to capture implicit affective states (Balahur et al., 2011, 2012; Udochukwu and He, 2015). Classical algorithms, such as Support Vector Machines (SVM) and Naive Bayes (NB), were combined with lexical features, syntactic patterns, and semantic resources to infer implicit emotions (Balahur et al., 2012; Riahi and Safari, 2016; Khoshnam and Baraani-Dastjerdi, 2022). These methods, although somewhat effective, have had difficulty with generalization due to the complexity of implicit emotional expressions and the absence of explicit emotional words.

The emergence of deep learning (DL) has introduced new avenues for recognizing implicit emotions in textual data, leveraging neural architectures such as Convolutional Neural Networks (CNNs) and Recurrent Neural Networks (RNNs) to capture complex linguistic and contextual patterns. Models such as LSTM and BiLSTM incorporated with attention mechanisms demonstrate improved performance by capturing temporal dependencies and contextual information (Rozental et al., 2018; Balazs et al., 2018; Chronopoulou et al., 2018; Rathnayaka et al., 2018; Zhou and Wu, 2018; Witon et al., 2018; Pecar et al., 2018; Fei et al., 2019). Recently, transformer-based models like BERT (Devlin et al., 2018) further enhanced the IER task by leveraging pre-trained embeddings and self-attention mechanisms, making them well suited for understanding implicit cues (Khoshnam et al., 2022; Qian et al., 2023; Boutouta et al., 2025).

Transformer-based LLMs, such as ChatGPT, have significantly expanded the possibilities of the NLP field, demonstrating remarkable performance across a wide range of tasks. These tasks include text understanding and generation (Mitrović et al., 2023; Gao et al., 2024), machine translation (Peng

et al., 2023), sentiment analysis (Buscemi and Proverbio, 2024), and semantic role labeling (Senator et al., 2025). Their strong generalization capabilities and the ability to capture contextual nuances enable more accurate emotion identification without requiring additional training (Kadiyala, 2024; Banimelhem and Amayreh, 2023; Lee et al., 2024; Hong et al., 2024; Liu et al., 2024). A recent study by Hong et al. (2024) introduced a method that addresses the complex and ambiguous nature of human emotions by using LLMs for ER. The approach considers multiple emotion labels and the intricate nature of emotional expressions. Another work proposed EmoLLMs (Liu et al., 2024), a series of open-source instruction-following LLMs fine-tuned for comprehensive affective analysis. These models are trained on a diverse dataset covering various classification and regression tasks related to emotions, enhancing their applicability in ER tasks. In another study (Wake et al., 2023), the performance of ChatGPT in the area of emotion detection was assessed on a variety of datasets, including IEMOCAP and DailyDialog. ChatGPT was able to classify text with emotional labels in both zero-shot and fine-tuning settings.

Despite these significant advances in ER, existing studies have predominantly focused on EER, with limited attention given to the IER task. To the best of our knowledge, none of the existing works have comprehensively addressed the unique challenges of IER, nor have they fully examined the potential of LLMs, such as ChatGPT, within this context.

## 3 Datasets

Two datasets were used: the WASSA-2018 Implicit Emotions Shared Task (IEST) dataset (Klinger et al., 2018) and the International Survey on Emotion Antecedents and Reactions (ISEAR) dataset (Scherer, 2005). Both datasets are widely used for ER, but differ significantly in terms of domain, format, and emotion expression. Table 5 in the Appendix A presents a brief comparison between the IEST and ISEAR datasets, while the label distributions are shown in Fig. 1 and Fig. 2.

### 3.1 IEST

The IEST dataset[1], introduced by Klinger et al. (2018), was developed for the WASSA-2018 Implicit Emotions Shared Task. It is a large automat-

ically labeled dataset of 191,731 English tweets, split into 153,600 for training, 9,600 for validation, and 28,800 for testing. Each tweet is annotated with one of Ekman's six basic emotions: anger, disgust, fear, joy, sadness, or surprise. Given computational constraints, only the testing set of the IEST dataset was used in this study.

To simulate implicit emotion scenarios, each tweet in the dataset has had its explicit emotion word masked and replaced with a placeholder token [#TARGETWORD#]. This design forces models to rely only on contextual cues to infer the underlying emotion, making it particularly suited for research on emotion understanding in indirect and implicit expressions. Some examples from the dataset are provided in Table 6 in the Appendix A.

### 3.2 ISEAR

The ISEAR dataset[2], introduced by Scherer (2005), is a manually labeled dataset collected as part of a psychological study aimed at exploring emotional experiences across cultures. The data were gathered from over 3,000 participants in 26 countries, all of whom had university-level education and were fluent in English. Each participant was asked to describe situations in which they had personally experienced one of seven emotions: joy, fear, anger, sadness, disgust, shame, and guilt. In total, the dataset contains approximately 7,666 instances, making it one of the most widely cited benchmarks for ER in psychology and affective computing tasks. Examples from the dataset are provided in Table 7 in Appendix A.

### 3.3 Data pre-processing

To align the ISEAR dataset with the IER task, we applied an additional filtering step: we ensured that none of the selected instances contained explicit emotion words. This pre-processing step allows us to reuse the ISEAR as a proxy dataset for IER, focusing only on instances where emotions must be inferred from the described context rather than directly stated. Furthermore, both datasets were subjected to standard pre-processing steps, including the removal of HTML tags, URLs, emojis, and extra spaces, as well as the correction of inconsistent punctuation.

---

[1]https://implicitemotions.wassa2018.com/data/

[2]https://github.com/sinmaniphel/py_isear_dataset

Figure 1: Distribution of emotion labels in the IEST dataset.



Figure 2: Distribution of emotion labels in the ISEAR dataset.

## 4 Experimental Methodology

We investigate two prominent approaches for text classification in the context of IER: (1) fine-tuning masked language models, and (2) prompt-based interaction with causal language LLMs. For the first approach, we employ BERT and RoBERTa, both pre-trained transformer encoders that are fine-tuned on task-specific data. These models have been widely recognized for their ability to capture contextual semantics and perform well across a range of NLP tasks. In this setup, the models are initialized with pre-trained weights and then fine-tuned using supervised learning on labeled emotion data. In contrast, the second approach utilizes Chat-GPT, a large, decoder-based LLM, accessed via zero-shot prompting. Rather than fine-tuning the model, we interact with ChatGPT using carefully crafted prompts that define the task and specify the desired output format. This method evaluates Chat-GPT's ability to generalize to the IER task without the need for additional training or fine-tuning.

By comparing these two paradigms, we aim to assess the trade-offs in performance, flexibility, and data efficiency when applied to implicit emotion classification.

### 4.1 Models

#### 4.1.1 BERT

A state-of-the-art NLP model introduced by Google in 2018 (Devlin et al., 2018) revolutionized the field by leveraging a bidirectional transformer architecture, which allows it to capture context from both the left and right of a word simultaneously. BERT is pre-trained on large text corpora using two key objectives: Masked Language Modeling (MLM), where it predicts randomly masked words within a sentence, and Next Sentence Prediction (NSP), where it learns to determine whether one sentence logically follows another. These pre-training tasks enable BERT to develop a deep understanding of both semantic meaning and syntactic structure in natural language.

#### 4.1.2 RoBERTa

Developed by Facebook AI in 2019 upon the foundational BERT architecture (Liu, 2019). RoBERTa improves and optimizes BERT's pre-training process by removing the NSP objective, training on significantly larger datasets, and employing dynamic masking during pre-training. These enhancements lead to improved performance on a wide range of natural language understanding tasks.

#### 4.1.3 ChatGPT

An advanced LLM developed by OpenAI in November 2022 (OpenAI), based on the Generative Pre-trained Transformer (GPT) architecture, a causal variant of the transformer neural network that has become the industry standard for a wide range of NLP tasks (Gillioz et al., 2020). Unlike masked language models, GPT models are trained in an autoregressive manner to predict the next token in a sequence, enabling strong generative and contextual reasoning abilities. ChatGPT was trained on a vast and diverse corpus, including academic texts, literary works, and large-scale web content, which equips it with broad linguistic and world knowledge. One of its key features is its ability to generate coherent, contextually relevant, and human-like responses to user input. Through interactive prompt-based querying, ChatGPT can adapt flexibly to new tasks without the need for additional fine-tuning.

### 4.2 Evaluation Approaches

#### 4.2.1 Fine-tuning encoder-based models

Fine-tuning involves adapting the pre-trained language models to a specific task by training them

on a smaller, task-specific dataset. This process requires significantly less data compared to training a model from scratch, thanks to the rich linguistic knowledge already encoded in the pre-trained model parameters. Regarding encoder-based models, we explored a range of hyperparameters configurations to optimize performance. Specifically, we experimented with different learning rates (1e-5, 2e-5, and 3e-5), batch sizes (16 and 32), training durations (ranging from 3 to 6 epochs), and maximum sequence lengths (64, 128, and 512). Additionally, we compared different model variants, including base and large versions, to assess their suitability for the IER task.

Each configuration was evaluated using a 10% development split of the training data, and the optimal setup was selected based on the macro F1-score. The chosen hyperparameters were validated across various random seeds to ensure robustness. Table 1 summarizes both the tested and optimal hyperparameter configurations.

We used the pre-trained "bert-base-uncased" and "roberta-base" models from the Huggingface Transformers library. The models consist of 12 transformer layers, a hidden size of 768, and 12 attention heads. For both models, we appended a dense layer with a softmax activation function for classification. The models were trained for 4 epochs using a batch size of 32 and a maximum sequence length of 128. Training was performed using the Adam optimizer and categorical cross-entropy loss. We evaluated these models on held-out test sets comprising 10% of the IEST and ISEAR datasets.

### 4.2.2 Prompt Design for Zero-Shot IER

For ChatGPT, we evaluated its zero-shot performance on test sets consisting of 600 and 700 instances from the IEST and ISEAR datasets, respectively (100 instances per emotion). As a proprietary model, ChatGPT was accessed via its chatbot interface using the GPT-4 Turbo version. To eliminate potential influence from prior context, each input was submitted in a separate chat session, ensuring full isolation between predictions.

Carefully designed zero-shot prompts are essential for enabling LLMs to generalize effectively across diverse domains (Team et al., 2023). We prompted ChatGPT with a text sample, a predefined list of emotion labels, task-specific instructions, and a set of output constraints. The prompts were iteratively designed and refined to align with the task's unique demands, namely, detecting emo-

tional states without the presence of explicit emotion words. Early versions of the prompt included basic task instructions. However, we observed improved performance when the implicit nature of the task was explicitly stated, when emotion label choices were clearly specified, and when the model's role was defined. For example, we experimented with formulations such as "the emotion is implied rather than stated." Additionally, we consistently framed the model as an "expert in implicit emotion recognition" at the beginning of each interaction to guide its behavior.

After multiple iterations, the final prompt adopted was:

---

**Role**: *You are an expert in implicit emotion recognition.*
**Prompt**: *The following sentence contains an emotion that is expressed implicitly. Based on context alone, identify the most likely emotion. Choose only one from: [Emotion List].*
*Respond with the emotion without any explanation.*
**Text**: *[example text]*

---

This final format was selected after testing several prompt versions on a development subset of the IEST and ISEAR datasets, evaluating performance manually and through agreement with gold-standard labels. We observed that prompting clarity, emotion list formatting, and explicit task framing significantly affected model responses.

### 4.3 Evaluation Metrics

Classification problem's performance is evaluated using a set of metrics. In our case, we use the accuracy and the macro average precision, recall, and F1-score. Each metric is defined in accordance with the following equations: (1), (2), (3), and (4), respectively. Where TP, TN, FP, and FN represent the number of True Positives, True Negatives, False Positives, and False Negatives, respectively.

$$Accuracy = \frac{TP + TN}{TP + TN + FP + FN} \quad (1)$$

$$Precision = \frac{TP}{TP + FP} \quad (2)$$

$$Recall = \frac{TP}{TP + FN} \quad (3)$$

$$F1 - score = 2 \cdot \frac{Precision * Recall}{Precision + Recall} \quad (4)$$

| Hyperparameter | Tested Values | Optimal Value |
|---|---|---|
| Learning Rate | 1e-5, 2e-5, 3e-5 | 1e-5 |
| Loss Function | Categorical Cross-Entropy | Categorical Cross-Entropy |
| Optimizer | Adam | Adam |
| Batch Size | 16, 32 | 32 |
| Epochs | 3, 4, 5, 6 | 4 |
| Max Length | 64, 128, 512 | 128 |

Table 1: Hyperparameter settings

## 5 Results and Discussion

Table 2 provides a concise overview of the performance of encoder-based LLMs (BERT and RoBERTa) and the decoder-based LLM (ChatGPT), using different approaches (fine-tuning and zero-shot prompting) across the IEST and ISEAR datasets.

### 5.1 Adaptation and Generalization

ChatGPT achieves the highest accuracy (77.14%) and F1-score (77.00%) on the ISEAR dataset, outperforming fine-tuned BERT and RoBERTa on the same dataset, which achieve an accuracy of 70.36% and 70.84%, respectively. However, on the IEST dataset, fine-tuned RoBERTa performs best (66.88% accuracy, 66.67% F1-score), while ChatGPT's performance drops significantly (54.17% accuracy, 54.93% F1-score). These results highlight a fundamental distinction between generalization and adaptation in IER. ChatGPT, as a causal language model, leverages broad pre-training to generalize well on datasets like ISEAR, where contextual cues align with its prior knowledge. In contrast, fine-tuned BERT and RoBERTa models, as masked language models, demonstrate superior adaptation to domain-specific constraints in IEST, where emotional keywords are masked, and cues are subtle. The masked modeling architecture, coupled with task-specific fine-tuning, equips these models with the ability to capture fine-grained contextual dependencies tailored to the dataset's structure, whereas ChatGPT's causal generation approach, optimized for predicting the next token, may be less effective in such constrained contexts. This performance gap underscores how model architecture and training paradigms interact with dataset characteristics to shape success in IER.

### 5.2 Performance Variation Across Datasets

As we show in Fig. 3, all models consistently performed better on the ISEAR dataset than on the IEST dataset. A possible reason for this finding is the contrast between the two datasets. While the ISEAR dataset was originally developed for general ER, we adapted it for the IER task by excluding any instances containing explicit emotion words (as noted in Section 3.3). This ensured that emotional states had to be inferred from contextual and situational cues rather than directly stated. Nevertheless, ISEAR dataset remains more clear, formal, consisting of well-structured, self-reported emotional experiences. These descriptions tend to be complete, coherent, and grammatically consistent. In contrast, the IEST dataset is derived from social media (tweets), which are often informal, fragmented, noisy, and contextually ambiguous. In addition, tweets may include slang, sarcasm, or cultural references that are not easily interpreted without broader context. This shift in genre presents additional challenges for IER, as models must not only infer unstated emotions but also navigate less structured and noisier linguistic input. We include representative examples from both datasets and a comparative table in Appendix A to illustrate these variations.

### 5.3 Emotional Implicitness

When considering the implicit emotional expression in each dataset, the IEST dataset represents masked emotion as a proxy for implicit emotion, where explicit emotion words were originally present in the sentence but have been deliberately removed. This deliberate omission weakens contextual support, forcing models to infer emotions from incomplete or ambiguous linguistic cues. In contrast, ISEAR contains naturally implicit emotions embedded within coherent, narrative-style descriptions of personal experiences. These richer and more structured contexts provide clearer situational signals, which both fine-tuned models and ChatGPT exploited to infer emotions more effectively. This distinction highlights how the availability and

| Approach | Model | ISEAR | | IEST | |
|---|---|---|---|---|---|
| | | Acc (%) | F1 (%) | Acc (%) | F1 (%) |
| *Fine-tuned* | BERT | 70.36 | 69.66 | 62.02 | 61.52 |
| | RoBERTa | 70.84 | 70.34 | **66.88** | **66.67** |
| *Zero shot* | ChatGPT | **77.14** | **77.00** | 54.17 | 54.93 |

Table 2: Performance comparison of different models on ISEAR and IEST datasets.



Figure 3: Performance comparison of different models on ISEAR and IEST datasets.

quality of contextual information directly shape the difficulty of implicit emotion recognition, with IEST posing a greater challenge due to its sparse and less informative cues.

## 5.4 Performance on Individual Emotions

Table 3 presents the performance of the three models for each emotion on the ISEAR dataset. As indicated, ChatGPT demonstrates superior performance for the majority of the emotions. Specifically, it achieves an F1-score of 94% for the emotion 'joy' and 84% for 'fear,' significantly outperforming the fine-tuned BERT and RoBERTa models. This demonstrates its strong ability to recognize emotions with clear contextual cues. However, all models show relatively poor performance on shame compared to other emotions, with F1-scores of 47%, 51%, and 59% for BERT, RoBERTa, and ChatGPT, respectively. The lower scores for shame reflect the challenge of detecting emotions that are highly implicit, underscoring the critical role of contextual clarity in IER.

The results on the IEST dataset are presented in Table 4, revealing a different trend compared to the ISEAR dataset. In this case, RoBERTa achieves the best performance across most emotions, with F1-scores of 77%, 78%, and 58% for the emotions 'joy', 'fear', and 'anger', respectively. These results significantly outperform those of the BERT and

ChatGPT models.

We also noticed significant variation when examining performance based on individual emotion labels. For example, in the zero-shot experiments on the ISEAR dataset, the recognition performance (F1-score) for 'joy' was around 94%, while it was below 60% for 'shame'. Similarly, on the IEST dataset, the F1-score for 'fear' was around 66%, while it was below 46% for 'anger'. In the fine-tuning approach, we observed that IER performance varied significantly across datasets, even for similar emotions. For instance, in the ISEAR dataset, the recognition performance (F1-score) for 'joy' was around 86% and 92% for BERT and RoBERTa models, respectively, while in the IEST dataset, the F1-scores for the same emotion (joy) were only around 70% and 77% for the same models, respectively. Notably, this tendency is also observed in the zero-shot condition with ChatGPT. For example, in the ISEAR dataset, the F1-score for 'anger' was around 73%, while in the IEST dataset, it was only around 45% with ChatGPT. This contrast in performance demonstrated that IER is a challenging task, with performance varying significantly depending on emotions, datasets, and models.

## 6 Conclusions and future works

In this study, we examined the effectiveness of two different architectures of LLMs for recognizing implicit emotions: masked language models, including BERT and RoBERTa, via a series of fine-tuning experiments, and causal language models, represented by ChatGPT, using a zero-shot prompting approach. The models were tested on two datasets: IEST and ISEAR. Both datasets are widely used for ER, but they differ significantly in terms of domain, format, and emotion expression. Our findings indicate that BERT-based fine-tuned models, particularly RoBERTa, excel at capturing implicit emotional cues. In contrast, zero-shot ChatGPT delivers promising results for certain emotion categories but struggles with more com-

| Model | Joy | | | Fear | | | Anger | | | Sadness | | | Disgust | | | Shame | | | Guilt | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | P | R | F1 | P | R | F1 | P | R | F1 | P | R | F1 | P | R | F1 | P | R | F1 | P | R | F1 |
| BERT | 0.82 | 0.91 | 0.86 | **0.78** | 0.78 | 0.78 | 0.63 | 0.57 | 0.60 | 0.70 | 0.83 | 0.76 | 0.76 | 0.76 | 0.76 | 0.65 | 0.37 | 0.47 | 0.58 | **0.72** | 0.64 |
| RoBERTa | 0.93 | 0.91 | 0.92 | 0.65 | 0.79 | 0.71 | 0.61 | 0.63 | 0.71 | 0.72 | 0.83 | 0.78 | 0.71 | **0.80** | 0.75 | 0.63 | 0.43 | 0.51 | **0.69** | 0.57 | **0.78** |
| ChatGPT | **0.96** | **0.92** | **0.94** | **0.78** | **0.90** | **0.84** | **0.72** | **0.73** | **0.73** | **0.74** | **0.92** | **0.82** | **0.91** | 0.70 | **0.79** | **0.67** | **0.53** | **0.59** | 0.65 | **0.72** | 0.68 |

Table 3: Performance comparison of different models per emotion on ISEAR dataset.

| Model | Joy | | | Fear | | | Anger | | | Sad | | | Disgust | | | Surprise | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | P | R | F1 | P | R | F1 | P | R | F1 | P | R | F1 | P | R | F1 | P | R | F1 |
| BERT | 0.66 | **0.75** | 0.70 | 0.64 | **0.80** | 0.71 | **0.57** | 0.47 | 0.51 | 0.60 | 0.53 | 0.57 | **0.67** | 0.53 | 0.57 | 0.57 | **0.64** | **0.60** |
| RoBERTa | **0.79** | **0.75** | **0.77** | **0.78** | 0.78 | **0.78** | 0.55 | **0.62** | **0.58** | **0.63** | **0.66** | **0.65** | 0.65 | **0.59** | **0.62** | 0.61 | 0.59 | **0.60** |
| ChatGPT | 0.65 | 0.61 | 0.63 | 0.77 | 0.57 | 0.66 | 0.36 | **0.62** | 0.45 | 0.48 | 0.53 | 0.50 | 0.60 | 0.52 | 0.56 | **0.66** | 0.40 | 0.50 |

Table 4: Performance comparison of different models on emotion classification on IEST dataset.

plex and context-dependent cases, where its performance declines noticeably. These results highlight the strengths of fine-tuned, medium-sized language models in handling IER tasks, while also underscoring the potential of zero-shot LLMs for emotions that are simpler or positively valenced. However, progress in IER remains constrained by the scarcity of high-quality datasets. Emotions are often conveyed indirectly, and building datasets that capture this nuance without relying on explicit markers is inherently challenging. Despite its limitations, the IEST dataset serves as a practical proxy by simulating implicitness through masked emotion words, offering a controlled evaluation setting.

Future research would benefit from the development of more diverse and realistic datasets for IER, as current resources are limited and often fail to capture the nuanced and context-dependent nature of implicit expressions. In addition to zero-shot prompting, we will explore alternative strategies such as few-shot learning and fine-tuning with LLMs, aiming to combine the adaptability of prompt-based approaches with the task-specific precision of supervised learning. Finally, addressing the persistent challenge of detecting socially complex emotions, such as *shame*, *guilt*, remains an important direction for future investigation, as these emotions often rely on subtle discourse cues and cultural context.

## Limitations

Despite the promising results presented in this study, some limitations should be considered. First, the evaluation was restricted to two benchmark datasets, ISEAR and IEST, which, although widely used in the field of ER, may not comprehensively reflect the variability of implicit emotional expres-

sions encountered in real-world scenarios, particularly in social media or multilingual contexts. This raises concerns regarding the generalization of the findings. Second, while LLMs exhibit an ability to capture certain contextual and cultural cues, their comprehension remains limited in the presence of more nuanced expressions such as sarcasm, idioms, or domain-specific references, which are common in implicit emotional content. Lastly, the lack of interpretability remains a critical challenge, particularly with generative models like ChatGPT, which operate as black boxes. This opacity hinders the ability to understand or explain model decisions, posing a barrier to trust and transparency in practical applications.

## References

Nourah Alswaidan and Mohamed El Bachir Menai. 2020. A survey of state-of-the-art approaches for emotion recognition in text. *Knowledge and Information Systems*, 62(8):2937–2987.

Alexandra Balahur, Jesús M Hermida, and Andrés Montoyo. 2011. Detecting implicit expressions of sentiment in text based on commonsense knowledge. In *Proceedings of the 2nd workshop on computational approaches to subjectivity and sentiment analysis (WASSA 2.011)*, pages 53–60.

Alexandra Balahur, Jesús M Hermida, and Andrés Montoyo. 2012. Detecting implicit expressions of emotion in text: A comparative analysis. *Decision support systems*, 53(4):742–753.

Jorge Balazs, Edison Marrese-Taylor, and Yutaka Matsuo. 2018. IIIDYT at IEST 2018: Implicit emotion classification with deep contextualized word representations. In *Proceedings of the 9th Workshop on Computational Approaches to Subjectivity, Sentiment and Social Media Analysis*, pages 50–56, Brussels, Belgium. Association for Computational Linguistics.

Omar Banimelhem and Wlla Amayreh. 2023. The performance of chatgpt in emotion classification. In *2023 14th International Conference on Information and Communication Systems (ICICS)*, pages 1–4. IEEE.

Carmen Bisogni, Lucia Cimmino, Maria De Marsico, Fei Hao, and Fabio Narducci. 2023. Emotion recognition at a distance: The robustness of machine learning based on hand-crafted facial features vs deep learning models. *Image and Vision Computing*, 136:104724.

Hanane Boutouta, Abdelaziz Lakhfif, Ferial Senator, and Chahrazed Mediani. 2025. A transformer-based hybrid model for implicit emotion recognition in arabic text. *Engineering, Technology & Applied Science Research*, 15(3):23834–23839.

Alessio Buscemi and Daniele Proverbio. 2024. Chatgpt vs gemini vs llama on multilingual sentiment analysis. *arXiv preprint arXiv:2402.01715*.

Alexandra Chronopoulou, Aikaterini Margatina, Christos Baziotis, and Alexandros Potamianos. 2018. NTUA-SLP at IEST 2018: Ensemble of neural transfer methods for implicit emotion classification. In *Proceedings of the 9th Workshop on Computational Approaches to Subjectivity, Sentiment and Social Media Analysis*, pages 57–64, Brussels, Belgium. Association for Computational Linguistics.

Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. 2018. Bert: Pre-training of deep bidirectional transformers for language understanding. *arXiv preprint arXiv:1810.04805*.

Bosheng Ding, Chengwei Qin, Linlin Liu, Yew Ken Chia, Shafiq Joty, Boyang Li, and Lidong Bing. 2022. Is gpt-3 a good data annotator? *arXiv preprint arXiv:2212.10450*.

Hao Fei, Yafeng Ren, and Donghong Ji. 2019. Implicit objective network for emotion detection. In *CCF International Conference on Natural Language Processing and Chinese Computing*, pages 647–659. Springer.

Ge Gao, Jongin Kim, Sejin Paik, Ekaterina Novozhilova, Yi Liu, Sarah T Bonna, Margrit Betke, and Derry Tanti Wijaya. 2024. Enhancing emotion prediction in news headlines: Insights from chatgpt and seq2seq models for free-text generation. In *Proceedings of the 2024 Joint International Conference on Computational Linguistics, Language Resources and Evaluation (LREC-COLING 2024)*, pages 5944–5955.

Anthony Gillioz, Jacky Casas, Elena Mugellini, and Omar Abou Khaled. 2020. Overview of the transformer-based models for nlp tasks. In *2020 15th Conference on computer science and information systems (FedCSIS)*, pages 179–183. IEEE.

Vishal Goar, Nagendra Singh Yadav, and Pallavi Singh Yadav. 2023. Conversational ai for natural language

processing: An review of chatgpt. *International Journal on Recent and Innovation Trends in Computing and Communication*, 11:109–117.

Amr Hendy, Mohamed Abdelrehim, Amr Sharaf, Vikas Raunak, Mohamed Gabr, Hitokazu Matsushita, Young Jin Kim, Mohamed Afify, and Hany Hassan Awadalla. 2023. How good are gpt models at machine translation? a comprehensive evaluation. *arXiv preprint arXiv:2302.09210*.

Xin Hong, Yuan Gong, Vidhyasaharan Sethu, and Ting Dang. 2024. Aer-llm: Ambiguity-aware emotion recognition leveraging large language models. *arXiv preprint arXiv:2409.18339*.

Ram Mohan Rao Kadiyala. 2024. Cross-lingual emotion detection through large language models. In *Proceedings of the 14th Workshop on Computational Approaches to Subjectivity, Sentiment, & Social Media Analysis*, pages 464–469.

Fereshteh Khoshnam and Ahmad Baraani-Dastjerdi. 2022. A dual framework for implicit and explicit emotion recognition: An ensemble of language models and computational linguistics. *Expert Systems with Applications*, 198:116686.

Fereshteh Khoshnam, Ahmad Baraani-Dastjerdi, and MJ Liaghatdar. 2022. Cefer: A four facets framework based on context and emotion embedded features for implicit and explicit emotion recognition. *arXiv preprint arXiv:2209.13999*.

Roman Klinger, Orphée De Clercq, Saif Mohammad, and Alexandra Balahur. 2018. IEST: WASSA-2018 implicit emotions shared task. In *Proceedings of the 9th Workshop on Computational Approaches to Subjectivity, Sentiment and Social Media Analysis*, pages 31–42, Brussels, Belgium. Association for Computational Linguistics.

Sheetal Kusal, Shruti Patil, Ketan Kotecha, Rajanikanth Aluvalu, and Vijayakumar Varadarajan. 2021. Ai based emotion detection for textual big data: Techniques and contribution. *Big Data and Cognitive Computing*, 5(3):43.

Sanghyub John Lee, Hyunseo Tony Lee, and Kiseong Lee. 2024. Enhancing emotion detection through chatgpt-augmented text transformation in social media text. In *2024 33rd IEEE International Conference on Robot and Human Interactive Communication (ROMAN)*, pages 872–879. IEEE.

Yinhan Liu. 2019. Roberta: A robustly optimized bert pretraining approach. *arXiv preprint arXiv:1907.11692*, 364.

Zhiwei Liu, Kailai Yang, Qianqian Xie, Tianlin Zhang, and Sophia Ananiadou. 2024. Emollms: A series of emotional large language models and annotation tools for comprehensive affective analysis. In *Proceedings of the 30th ACM SIGKDD Conference on Knowledge Discovery and Data Mining*, pages 5487–5496.

Sandra Mitrović, Davide Andreoletti, and Omran Ayoub. 2023. Chatgpt or human? detect and explain. explaining decisions of machine learning model for detecting short chatgpt-generated text. *arXiv preprint arXiv:2301.13852*.

OpenAI. Chatgpt. Accessed: 2025-02-23.

Udochukwu Orizu. 2018. *Implicit emotion detection in text*. Ph.D. thesis, Aston University.

Samuel Pecar, Michal Farkas, Marián Šimko, Peter Lacko, and Maria Bielikova. 2018. Nl-fiit at iest-2018: Emotion recognition utilizing neural networks and multi-level preprocessing. In *Proceedings of the 9th Workshop on Computational Approaches to Subjectivity, Sentiment and Social Media Analysis*, pages 217–223.

Keqin Peng, Liang Ding, Qihuang Zhong, Li Shen, Xuebo Liu, Min Zhang, Yuanxin Ouyang, and Dacheng Tao. 2023. Towards making the most of chatgpt for machine translation. *arXiv preprint arXiv:2303.13780*.

Yanjun Qian, Jin Wang, Dawei Li, and Xuejie Zhang. 2023. Interactive capsule network for implicit sentiment analysis. *Applied Intelligence*, 53(3):3109–3123.

Prabod Rathnayaka, Supun Abeysinghe, Chamod Samarajeewa, Isura Manchanayake, and Malaka Walpola. 2018. Sentylic at iest 2018: Gated recurrent neural network and capsule network-based approach for implicit emotion detection. *arXiv preprint arXiv:1809.01452*.

Nooshin Riahi and Pegah Safari. 2016. Implicit emotion detection from text with information fusion. *Journal of Advances in Computer Research*, 7(2):85–99.

Alon Rozental, Daniel Fleischer, and Zohar Kelrich. 2018. Amobee at iest 2018: Transfer learning from language models. *Proceedings of the 9th Workshop on Computational Approaches to Subjectivity, Sentiment and Social Media Analysis*.

Klaus R Scherer. 2005. What are emotions? and how can they be measured? *Social science information*, 44(4):695–729.

Ferial Senator, Abdelaziz Lakhfif, Imene Zenbout, Hanane Boutouta, and Chahrazed Mediani. 2025. Leveraging chatgpt for enhancing arabic nlp: Application for semantic role labeling and cross-lingual annotation projection. *IEEE Access*, 13:3707–3725.

Gemini Team, Rohan Anil, Sebastian Borgeaud, Jean-Baptiste Alayrac, Jiahui Yu, Radu Soricut, Johan Schalkwyk, Andrew M Dai, Anja Hauth, Katie Millican, et al. 2023. Gemini: a family of highly capable multimodal models. *arXiv preprint arXiv:2312.11805*.

Orizu Udochukwu and Yulan He. 2015. A rule-based approach to implicit emotion detection in text. In *Natural Language Processing and Information Systems: 20th International Conference on Applications of Natural Language to Information Systems, NLDB 2015, Passau, Germany, June 17-19, 2015, Proceedings 20*, pages 197–203. Springer.

Naoki Wake, Atsushi Kanehira, Kazuhiro Sasabuchi, Jun Takamatsu, and Katsushi Ikeuchi. 2023. Bias in emotion recognition with chatgpt. *arXiv preprint arXiv:2310.11753*.

Wojciech Witon, Pierre Colombo, Ashutosh Modi, and Mubbasir Kapadia. 2018. Disney at iest 2018: Predicting emotions using an ensemble. In *Proceedings of the 9th Workshop on Computational Approaches to Subjectivity, Sentiment and Social Media Analysis*, pages 248–253.

Qimin Zhou and Hao Wu. 2018. Nlp at iest 2018: Bilstm-attention and lstm-attention via soft voting in emotion classification. In *Proceedings of the 9th workshop on computational approaches to subjectivity, sentiment and social media analysis*, pages 189–194.

Li'an Zhu, Junjie Peng, and Huiran Zhang. 2025. Text-based sarcasm detection with emoji contradictory clues assisting. In *Advanced Intelligent Computing Technology and Applications*, pages 248–260, Singapore. Springer Nature Singapore.

# A Appendix

| Aspect | IEST | ISEAR |
|---|---|---|
| Type of Emotion | Implicit (emotion word masked) | Explicit (emotion word present) |
| Annotation | Automatically labeled | Manually labeled |
| Emotion Labels | Ekman's six basic emotions: anger, disgust, fear, joy, sadness, surprise | anger, disgust, fear, joy, sadness, shame, guilt |
| Genre | social media (Twitter) | Survey responses |
| Style | Informal, noisy, fragmented sentence | formal, structured, complete sentences |
| Text Length | Short (tweets, < 280 characters) | Medium (1–3 sentences per instance) |
| Size | 191,731 instances | 7,666 instances |
| Purpose | IER | ER |
| Contextual Clues | Sparse; relies on social and situational context | Rich descriptions of emotional experiences |

Table 5: Comparison between the IEST and ISEAR datasets.

| Emotion | Tweet |
|---|---|
| Anger | I get impatient and [#TARGETWORD#] when I'm hungry. |
| Disgust | So many people looked at me just [#TARGETWORD#] when I said that mustaches are hot. |
| Fear | So [#TARGETWORD#] that I'm not good enough |
| Joy | you're gonna be [#TARGETWORD#] when you realize you deserve to be. |
| Sadness | Very [#TARGETWORD#] when he goes on these tirades |
| Surprise | They just jealous, they get [#TARGETWORD#] when she pull up. |

Table 6: Example tweets from the IEST dataset, with the emotion word masked as [#TARGETWORD#].

| Emotion | Example |
|---|---|
| Joy | An encounter with a man whom I love, after a very long separation. |
| Fear | After mischievously ringing on the chemist's trade-entrance doorbell and getting caught by him. |
| Anger | At my Summer job, nobody looked after me in particular and I had to learn all on my own. |
| Sadness | After I had lived with my boyfriend in a foreign country for half a year, I saw that it was impossible for me to stay with him (for economic reasons). We separated although I loved him. |
| Disgust | A mother who shouts at her child for nothing. |
| Shame | During carnaval I danced for a few minutes normally I don't dance because I am rigid in my moving around during a dance, I stopped very soon. |
| Guilt | I speak harshly to my parents though they only mean my own good. |

Table 7: Example from the ISEAR dataset for each emotion label.

# Cross-Lingual Sentence-Level Skill Identification
# in English and Danish Job Advertisements

**Nurlan Musazade**
Åbo Akademi University
Turku, Finland
nurlan.musazade@abo.fi

**Mike Zhang**
Aalborg University
Copenhagen, Denmark
jjz@cs.aau.dk

**József Mezei**
Åbo Akademi University
Turku, Finland
Jozsef.Mezei@abo.fi

## Abstract

The increasing influence of artificial intelligence (AI), the availability of textual data, and large language models (LLMs) over the past decade is evident in the growth of scholarly work on identifying skills from job advertisements. In this work, we examine the detection of sentences that express skills as well as the explainability of model decisions with respect to their dependence on skill related tokens. We compare traditional machine learning (ML) approaches with a pretrained multilingual model and domain-adapted models for the task of English skill identification, and we assess the role of skill tokens in the classification process. We also investigate the ability of these models to generalize from English (EN) to Danish (DA) in both few-shot and zero-shot settings. Our findings indicate that both models achieve high performance in sentence classification achieving an $F_1$-score of 94% for EN and overall accuracy between 93%–94% for both EN and DA. The results show that traditional ML methods can remain relevant under certain circumstances reinforcing the importance of realistic baselines in the context of skill identification.

## 1 Introduction

With the technological advancements and labor market disruptions, the importance of identifying skill requirements in job advertisements rises for both job seekers and educational institutions (Brasse, 2024). Job advertisements are a critical source to study skill requirements (Khaouja et al., 2021; Senger et al., 2024; Zhang et al., 2022a). The most straightforward method is sentence-level skill identification (SI), where the task is to predict whether a sentence contains a skill or not (Khaouja et al., 2021).

Although SI methods have been studied for the English language (Tamburri et al., 2020; Khaouja et al., 2021; Leon et al., 2024a; Rosenberger, 2025), it is unclear how well this extends to other languages. We hypothesize, considering that some

hard skills (e.g., Python, Java) are defined the same across languages, that there is a high generalizability level of English-based models. We extend prior research by testing the capabilities of (domain-adapted) multilingual language models (Chung et al., 2021; Zhang et al., 2023) on the task of SI, exploring the cross-lingual generalization of both a simple statistical baseline and fine-tuned models on English, and analyzing the factors that influence model decisions when classifying sentences. We seek to answer the following research questions:

**RQ1** How effective are logistic regression, domain trained and pre-trained LLMs in classifying skill-related sentences in job advertisements?

**RQ2** To what extent do English-based sentence classification models perform on Danish skill identification in zero- and few-shot settings?

**RQ3** What is the contribution of skill tokens in the classification of skill-related sentences in job advertisements?

**Contributions.** In this work, we contribute the following by answering the RQs by showing that: (i) both the baseline and multilingual LMs show high performance on in-language (English) skill identification and cross-lingual skill identification for Danish, even with little training data; (ii) sentence-level SI models' reliance on skill tokens are low, highlighting the need to assess context-dependent trustworthiness and robustness of such approach in real-world scenarios.

## 2 Related Work

### 2.1 Classification-Based Skill Identification

Khaouja et al. (2021) identified four primary skill extraction techniques: skill counting, topic modeling, embeddings (skill or word-based), and ma-

| Class | Train | Dev. | Test |
|---|---|---|---|
| 0 - No skill | 6,753 | 2,634 | 2,699 |
| 1 - Skill | 10,421 | 3,359 | 3,056 |

Table 1: **Class Distribution.** Distribution of skill/no-skill sentences in the combined English dataset.

chine learning methods. From an ML perspective, there are two principal methods: Named Entity Recognition (NER), which classifies and extracts entities into predefined categories and content-based text classification (i.e., classifying whether a sentence contains a skill or not). We adopt the latter, applying sentence-level binary classification. Using sentence-level granularity in skill identification helps preserve word relationships and context.

Lin et al. (2023) and Rosenberger (2025) highlighted the critical role of removing irrelevant information for achieving high model performance. Lin et al. (2023) explored synthetic data generation to enhance data relevance. Rosenberger (2025) used classification approaches to filter noise from job ads, significantly improving recommendation accuracy. Their jobGBERT model achieved high accuracy (0.96-0.97 $F_1$ score) by truncating irrelevant content at the paragraph level. In contrast, our research addresses sentence-level classification, potentially enhancing model training simplicity and annotation efficiency. Additionally, while Rosenberger (2025) targeted the German language, our research investigates English and Danish datasets, exploring cross-lingual transfer and zero- or few-shot learning scenarios.

Leon et al. (2024b) conducted binary and multi-label classification using English and multilingual models on job ads. Facing imbalanced data with a prevalence of skill-absent sentences, the authors applied augmentation techniques. The accuracy of these models ranged from 94%–99%. They noted limitations in data availability and domain adaptation challenges and emphasized the need for further exploration of multilingual applicability and model explainability.

More recently, there has been increasing focus on computational efficiency in skill extraction from job advertisements (Sun et al., 2025; Vásquez-Rodríguez et al., 2024). For example, Vásquez-Rodríguez et al. (2024) compared various methods to distinguish their effectiveness and efficiency.

Our study builds upon these works, such as multilingual capabilities and explainability to improve

| Model | Class | F1 | P | R | Acc. |
|---|---|---|---|---|---|
| LR + TF-IDF | 0 | 0.89 | 0.91 | 0.87 | 0.90 |
| | 1 | 0.90 | 0.89 | 0.92 | |
| RemBERT | 0 | 0.93 | 0.96 | 0.90 | 0.93 |
| | 1 | 0.94 | 0.92 | 0.96 | |
| ESCOXLM-R | 0 | 0.93 | 0.95 | 0.91 | **0.94** |
| | 1 | 0.94 | 0.92 | 0.96 | |

Table 2: **English Results.** Performance of SI on the English test set in terms of $F_1$-score, precision (P), recall (R), and Accuracy (Acc.).

real-time skill extraction and job recommendation systems. For a more detailed survey of SI, we refer to Khaouja et al. (2021).

## 3 Methodology

Skill identification at the sentence level can be formulated as a binary classification task. Given a set of sentences $S = \{s_1, s_2, \ldots, s_n\}$, each sentence $s_i$ is associated with a binary label $y_i$:

$$
y_i = \begin{cases} 1 & \text{if sentence } s_i \text{ contains a skill mention,} \\ 0 & \text{otherwise.} \end{cases}
$$

The goal of this task is to train a classification model $f$ that accurately predicts the label $y_i$ given a sentence $s_i$: $f(s_i) \to y_i$. The objective is to accurately classify each sentence on whether it contains a skill or not.

### 3.1 Data

We used three real-world English datasets for training and evaluation: SkillSpan (Zhang et al., 2022a), Green (Green et al., 2022), and Sayfullina (Sayfullina et al., 2018). Each dataset includes separate training, development, and test splits. In terms of dataset size, Green has 9,968, Sayfullina 7,411, and SkillSpan 11,543 sentences. In SkillSpan, there are both skill and knowledge annotations and were merged into one positive class. For Green, only positive skill annotations are taken.

Sentences were reconstructed from the word-level tokenized lists and each was labeled as either containing a skill (1) or not (0). Table 1 shows the distribution across all English splits. The combined data was shuffled before training. For Danish, we used the Kompetencer dataset (Zhang et al., 2022b), which follows the same format. It includes 778 training, 346 validation, and 262 test sentences.

| Training | Class | RemBERT | | | | ESCOXLM-R | | | |
|---|---|---|---|---|---|---|---|---|---|
| | | F1 | P | R | Acc. | F1 | P | R | Acc. |
| DA | 0 | 0.89 | 0.88 | 0.89 | 0.82 | 0.89 | 0.79 | 1.00 | 0.79 |
| | 1 | 0.55 | 0.57 | 0.54 | | 0.00 | 0.00 | 0.00 | |
| EN | 0 | 0.95 | 0.98 | 0.92 | 0.92 | 0.95 | 0.97 | 0.92 | 0.92 |
| | 1 | 0.83 | 0.75 | 0.93 | | 0.82 | 0.75 | 0.91 | |
| EN + DA | 0 | 0.95 | 0.96 | 0.95 | 0.93 | 0.95 | 0.96 | 0.95 | 0.93 |
| | 1 | 0.83 | 0.81 | 0.85 | | 0.83 | 0.81 | 0.85 | |

Table 3: **Cross-lingual Results.** Cross-lingual performance on the Danish test set, broken down by training data, model and class

## 3.2 Models

A random baseline for the English test set yields around 0.51 accuracy. For a second baseline, we trained a TF-IDF-based logistic regression model with unigrams and bigrams (5,000 max features).

To conduct SI with language models, we experimented with three models in this study; namely RemBERT (Chung et al., 2021), ESCOXLM-R (Zhang et al., 2023) which is a domain-adapted XLM-R-based model (Conneau et al., 2020).

## 3.3 Training and Evaluation

For training the multilingual language models, we explored a limited range of hyperparameters and finalized on a learning rate of $1 \times 10^{-6}$, batch size of 16, weight decay on 0.01, and ten epochs with patience of 2. For both the English and few-shot Danish experiments, we kept these hyperparameters. We measured performance by $F_1$ score, precision, recall, and accuracy.

## 3.4 Explainability

For analysis, we investigate model explainability and compare model behavior. We use Integrated Gradients (IG; Sundararajan et al., 2017), implemented with the Captum library (Kokhlikyan et al., 2020), to analyze true positive (TP) predictions. The aim is to interpret which input tokens contribute most to skill predictions, improving transparency. IG satisfies two key axioms: Sensitivity, where a differing feature between inputs with different outputs must receive non-zero attribution; and Implementation Invariance, where models with identical outputs for all inputs yield the same attributions.

## 4 Results and Discussion

### 4.1 English Results

In Table 2, we show the main results for English of both the baselines and language models. We observe that all three models achieve a high F-score (0.89–0.94) as well as a high accuracy (0.90–0.94). In particular, from a cost–efficiency perspective, the baseline took a few minutes to train, whereas ESCOXLM-r and RemBERT took between 80–140 minutes, with equal performance.

Among the large language models, performance differences are minor but important, especially for correctly identifying skill-containing sentences. Both models perform equally well on the positive class. ESCOXLM-R achieves slightly higher overall accuracy and better recall for the negative class, while RemBERT shows marginally better precision on that class. Given the small differences and limited model tuning, no definitive conclusion can be drawn about one model being superior.

### 4.2 Danish Results

For our cross-lingual experiments, we have three setups

- **(DA)**: Fine-tuning the language models on the few-hundred Danish instances.
- **(EN)**: Fine-tuning the language models on English only and then apply it to the Danish test set.
- **(EN + DA)**: Fine-tuning both models on English and Danish and apply it to Danish.

In Table 3, we show that DA and EN demonstrate good performance; precision starts from 0.75 and minimum recall is 0.85 for the positive class. The ESCOXLM-R model failed to predicted no skill sentences with DA, whereas the RemBERT model's performance without the English language training is around 0.55 for the positive class.

| Model | Precision | | | | Recall set | F1 set |
|---|---|---|---|---|---|---|
| | Top-1 | Top-2 | Top-3 | Set | | |
| LR + TF-IDF | 0.485 | 0.278 | 0.255 | 0.225 | 0.808 | 0.303 |
| ESCOXLM-R | 0.101 | 0.172 | 0.213 | 0.234 | 0.576 | 0.284 |
| RemBERT | 0.098 | 0.170 | 0.209 | 0.231 | 0.576 | 0.281 |

Table 4: **Explainability.** Comparison of the matching True Skill with the tokens contributing positively for the TP

| Text | True Skills | Top Attributed Tokens |
|---|---|---|
| javascript reactjs java | [javascript, reactjs, java] | [javascript, reactjs, java] |
| - and yaml | [yaml] | [-, ya] |
| Strong knowledge of application data and infrastructure architecture disciplines | [application, data, and, infrastructure, architecture] | [•, knowledge, application, data, infrastructure, s] |
| Demonstrated experience of performing DevOps for platforms | [DevOps, for, platforms] | [•ted, experience, of, performing, Dev, for] |
| You are proficient in Python and English | [Python, English] | [You, profi, in, English] |

Table 5: **Dataset Sample.** Sample data rows.

Performance of the two models are similar in the other settings (i.e., EN and EN + DA). Interestingly, for the positive class, we see when the models have been fine-tuned on English data only (EN) it outperforms the DA setting (0.06–0.08 $F_1$ higher), likely due to more training data being available, indicating successful cross-lingual transfer.

### 4.3 Explainability

For analysis, we compare the contribution of actual skills to class prediction focusing on the TP class. Table 4 shows that the sets of words driving positive predictions are very similar across the two models and that ESCOXLM-R performance is only slightly higher. In the base model we multiply the TF-IDF values by the coefficients to measure each word's contribution to TP predictions and then compare that with word attribution in the LMs. Precision in the base model is higher for the top two tokens but it declines when we consider the full token set, ending up below the LMs. The $F_1$-score stays similar in the base model due to the high recall.

All models have a precision around 0.22–0.23 and a 0.58 recall for the LMs and 0.81 for the baseline. This suggests that classification does not depend mainly on semantic content or even on the explicit presence of a skill term in the sentence. These results align with the observation that ESCOXLM-R, despite its domain training, performs similarly to RemBERT and that logistic regression with TF-IDF narrows most of the gap with the LMs.

## 5 Conclusion

In this work, show the effectiveness of language models, including multilingual ones and a basic supervised ML model with TF-IDF, for the task of skill identification. All our models achieve around 90%-94% accuracy on English, with ESCOXLM-r as best-performing, indicating that the task is straightforward.

The performance of the supervised baseline model demonstrates that traditional approaches should still be considered after the introduction of the advanced architectures and models, and can be beneficial, for instance in the low resource settings. Furthermore, we show the effectiveness of multilingual LMs for cross-lingual transfer. We show that the performance of multilingual LMs is still high (92%), even though we only train on English data. As expected, there is overlap between skills between languages. This could particularly benefit low resource languages either in zero-shot setting or with the minimal training data.

In our explainability analysis, we show that the contribution of the skill tokens do not contribute that much to the actual correct prediction, warranting further investigation. We conducted analysis at the token level without cleaning data from special characters and stopwords, which poses a limitation in the evaluation of the explainability, while representing a realistic inputs (see Table 5). Computing attribution scores for each word, e.g., by summing tokens' attribution scores, may present a more reli-

able and interpretable definitive measurement.

For future research, we will consider more languages to investigate whether transfer still holds. Additionally, we considered all skills as one, but can also distinguish between hard and soft skills.

## Ethics Statement

For the identification of specific occupational skills in sentences, we do not foresee any ethical issues.

## References

Julia Brasse. 2024. Identification of future skills using data-driven methods: A systematic literature review and directions for future research. *Proceedings of the 57th Hawaii International Conference on System Sciences*.

Hyung Won Chung, Thibault Févry, Henry Tsai, Melvin Johnson, and Sebastian Ruder. 2021. Rethinking embedding coupling in pre-trained language models. In *9th International Conference on Learning Representations, ICLR 2021, Virtual Event, Austria, May 3-7, 2021*. OpenReview.net.

Alexis Conneau, Kartikay Khandelwal, Naman Goyal, Vishrav Chaudhary, Guillaume Wenzek, Francisco Guzmán, Edouard Grave, Myle Ott, Luke Zettlemoyer, and Veselin Stoyanov. 2020. Unsupervised cross-lingual representation learning at scale. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 8440–8451, Online. Association for Computational Linguistics.

Thomas Green, Diana Maynard, and Chenghua Lin. 2022. Development of a benchmark corpus to support entity recognition in job descriptions. In *Proceedings of the Thirteenth Language Resources and Evaluation Conference*, pages 1201–1208, Marseille, France. European Language Resources Association.

Imane Khaouja, Ismail Kassou, and Mounir Ghogho. 2021. A survey on skill identification from online job ads. *IEEE Access*, 9:118134–118153.

Narine Kokhlikyan, Vivek Miglani, Miguel Martin, Edward Wang, Bilal Alsallakh, Jonathan Reynolds, Alexander Melnikov, Natalia Kliushkina, Carlos Araya, Siqi Yan, and Orion Reblitz-Richardson. 2020. Captum: A unified and generic model interpretability library for pytorch.

Florin Leon, Marius Gavrilescu, Sabina-Adriana Floria, and Alina Adriana Minea. 2024a. Hierarchical classification of transversal skills in job advertisements based on sentence embeddings. *Information*, 15(3).

Florin Leon, Marius Gavrilescu, Sabina-Adriana Floria, and Alina Adriana Minea. 2024b. Hierarchical classification of transversal skills in job advertisements based on sentence embeddings. *Information*, 15(3):151.

Shiyong Lin, Yiping Yuan, Carol Jin, and Yi Pan. 2023. Skill graph construction from semantic understanding. In *Companion Proceedings of the ACM Web Conference 2023*, pages 978–982.

Rosenberger. 2025. Careerbert: Matching resumes to esco jobs in a shared embedding space for generic job recommendations. *Expert Systems with Applications*.

Luiza Sayfullina, Eric Malmi, and Juho Kannala. 2018. Learning representations for soft skill matching. In *Analysis of Images, Social Networks and Texts: 7th International Conference, AIST 2018, Moscow, Russia, July 5–7, 2018, Revised Selected Papers 7*, pages 141–152. Springer.

Elena Senger, Mike Zhang, Rob van der Goot, and Barbara Plank. 2024. Deep learning-based computational job market analysis: A survey on skill extraction and classification from job postings. In *Proceedings of the First Workshop on Natural Language Processing for Human Resources (NLP4HR 2024)*, pages 1–15, St. Julian's, Malta. Association for Computational Linguistics.

Ying Sun, Yang Ji, Hengshu Zhu, Fuzhen Zhuang, Qing He, and Hui Xiong. 2025. Market-aware long-term job skill recommendation with explainable deep reinforcement learning. *ACM Transactions on Information Systems*, 43(2):1–35.

Mukund Sundararajan, Ankur Taly, and Qiqi Yan. 2017. Axiomatic attribution for deep networks. In *Proceedings of the 34th International Conference on Machine Learning, ICML 2017, Sydney, NSW, Australia, 6-11 August 2017*, volume 70 of *Proceedings of Machine Learning Research*, pages 3319–3328. PMLR.

Damian A Tamburri, Willem-Jan Van Den Heuvel, and Martin Garriga. 2020. Dataops for societal intelligence: a data pipeline for labor market skills extraction and matching. In *2020 IEEE 21st International Conference on Information Reuse and Integration for Data Science (IRI)*, pages 391–394. IEEE.

Laura Vásquez-Rodríguez, Samuel Michel Bertrand Audrin, Samuele Galli, Julneth Rogenhofer, Jacopo Negro Cusa, and Lonneke van der Plas. 2024. Hardware-effective approaches for skill extraction in job offers and resumes. *RecSys in HR'24: The 4th Workshop on Recommender Systems for Human Resources, in conjunction with the 18th ACM Conference on Recommender Systems*.

Mike Zhang, Kristian Jensen, Sif Sonniks, and Barbara Plank. 2022a. SkillSpan: Hard and soft skill extraction from English job postings. In *Proceedings of the 2022 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pages 4962–4984, Seattle, United States. Association for Computational Linguistics.

Mike Zhang, Kristian Nørgaard Jensen, and Barbara Plank. 2022b. Kompetencer: Fine-grained skill classification in Danish job postings via distant supervision and transfer learning. In *Proceedings of the*

*Thirteenth Language Resources and Evaluation Conference*, pages 436–447, Marseille, France. European Language Resources Association.

Mike Zhang, Rob van der Goot, and Barbara Plank. 2023. ESCOXLM-R: Multilingual taxonomy-driven pre-training for the job market domain. In *Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 11871–11890, Toronto, Canada. Association for Computational Linguistics.

# Demographics and Democracy: Benchmarking LLMs' Gender Bias and Political Leaning in European Parliament

**Jinrui Yang**[*]    **Xudong Han**[†]    **Timothy Baldwin**[*†]

[*]School of Computing & Information Systems, The University of Melbourne

[†]Mohamed bin Zayed University of Artificial Intelligence, UAE

jinrui.yang@student.unimelb.edu.au    xudong.han@mbzuai.ac.ae    tb@ldwin.net

## Abstract

We introduce EuroParlVote, a novel benchmark for evaluating large language models (LLMs) in politically sensitive contexts. It links European Parliament debate speeches to roll-call vote outcomes and includes rich demographic metadata for each Member of the European Parliament (MEP), such as gender, age, country, and political group. Using EuroParlVote, we evaluate state-of-the-art LLMs on two tasks—gender classification and vote prediction—revealing consistent patterns of bias. We find that LLMs frequently misclassify female MEPs as male and demonstrate reduced accuracy when simulating votes for female speakers. Politically, LLMs tend to favor centrist groups while underperforming on both far-left and far-right ones. Proprietary models like GPT-4o outperform open-weight alternatives in terms of both robustness and fairness. We release the EuroParlVote dataset, code, and demo to support future research on fairness and accountability in NLP within political contexts.

## 1 Introduction

With growing interest in applying natural language processing (NLP) methods to political discourse, recent studies have revealed persistent gender bias in the European Parliament. During parliamentary debates, certain subgroups — such as women, junior members, and representatives from smaller member states — receive disproportionately less attention and visibility (Walter et al., 2023). Similarly, gender bias has been shown to persist in political news coverage, with systematic disparities in word choice, sentiment, and framing across ideological lines, even when explicit gender markers are removed (Davis et al., 2022).

Meanwhile, recent studies have highlighted that many NLP technologies, including large language models (LLMs), exhibit measurable political biases, often leaning towards left-liberal viewpoints in their responses to political discourse (Rozado,

2024a; Feng et al., 2023b; Santurkar et al., 2024). However, these findings predominantly focus on U.S.-centric political contexts, for instance, Potter et al. (2024) analyzes discourse surrounding the 2024 U.S. presidential election. In contrast, this paper shifts the focus to the European Parliament, where we investigate how LLMs interpret and predict political behavior in a multilingual, multi-party democratic setting. We are interested in whether gender and ideological bias patterns observed in U.S. political contexts similarly manifest in the European setting.

Our study explores this question by introducing a novel EU voting dataset that links roll-call votes with corresponding debate speeches and detailed demographic information of each Member of European Parliament (MEP). We benchmark several LLMs on two tasks: predicting the gender of MEPs based on their speech, and simulating MEP voting behavior from debate content.

First, in Section 3, we construct a multilingual benchmark — covering 24 official EU languages — that links 22K European Parliament debate speeches to 969 corresponding roll-call votes. We further enrich the dataset with annotations about MEPs, including gender (male/female);[1] political group (across 8 groups including nonattached members); age (ranging from 25 to 83); and country (from 27 EU member states and one former member state, United Kingdom), enabling demographically-aware political modeling.

In Section 4, we analyze gender bias in LLMs (GPT-4o, LLaMA-3.2-3B, Claude-3.5, Gemini-2.5-Flash, and Mistral-large) within the context of the European Parliament. Specifically, we conduct two experiments: first, we ask LLMs to predict the gender of MEPs based solely on their debate speeches; and second, we provide the debate speeches, top-

---

[1]We acknowledge gender is non-binary, but use a male/female classification here, as it is an accurate representation of past and present MEPs.

ics, and the MEP's gender, and ask the models to predict their voting behavior.

Our findings reveal a consistent male-biased pattern in LLMs: (1) female MEPs are disproportionately misclassified as male in the gender prediction task; and (2) when all MEPs are hypothetically assigned the gender "female", the voting prediction accuracy drops to its lowest, whereas assigning all MEPs the gender "male" yields the highest accuracy; and (3) proprietary LLMs, such as GPT-4o and Gemini-2.5, inherently exhibit lower gender misclassification rates compared to open-weight models like LLaMA-3.2. This highlights how gender assumptions implicitly embedded in LLMs can influence both demographic classification and downstream political prediction tasks.

We also implemented LoRA-based fine-tuning (Hu et al., 2021) on open-weight LLMs using annotated training examples to evaluate its impact on gender bias mitigation. However, our results indicate that LoRA does not reduce gender bias in either LLaMA-3.2 or Mistral-large. This observation aligns with findings from Ding et al. (2024), which report that LoRA does not exhibit a consistent pattern of amplifying or mitigating disparate impacts across demographic subgroups.

In Section 5, we investigate the political leanings of LLMs through two experiments. First, given the debate topic and the speech content, we prompt the models to simulate a vote as if they were the MEP delivering the speech. Second, we additionally provide the models with the MEP's political group information and prompt them to vote again.

Our findings indicate that all LLMs exhibit a left–centrist bias, as evidenced by: (1) higher voting prediction accuracy for left-leaning and centrist political groups; but (2) for ideologically extreme groups, far-right parties are simulated more accurately than far-left ones; ; and (3) similar to gender bias, open-weight LLMs exhibit more pronounced political bias compared to proprietary models.

We further evaluate an instruction-tuned setting in which political group identifiers are included in the input prompt. This setup improves prediction accuracy for ideologically extreme groups — particularly far-left and far-right parties — suggesting that explicit political context helps mitigate performance disparities across ideological lines.

To the best of our knowledge, this is the first study to systematically benchmark both gender bias and political leaning in LLMs within the context of the European Parliament. Our results under-score the complexity of assessing and mitigating fairness concerns in LLM methods. We release EuroParlVote[2] and code[3] to support future research on fairness, transparency, and robustness in political NLP.

## 2 Related Work

### 2.1 Gender Bias of LLMs

Recent research has highlighted that LLMs often perpetuate, and sometimes amplify, gender stereotypes and biases. Kotek et al. (2023) propose a novel testing paradigm designed to probe for gender bias using linguistic constructions unlikely to be explicitly present in training data. Their study finds that LLMs frequently rely on gender stereotypes in completing tasks and that their justifications often cite faulty reasoning or make explicit reference to the stereotypes themselves. This suggests that even state-of-the-art LLMs, despite advancements enabled by techniques such as reinforcement learning with human feedback (RLHF) (Christiano et al., 2017), still encode and reproduce biased social patterns present in their training data. The authors argue that such biases reflect the "collective intelligence" of Western society as captured in large-scale text data, and call for improved diagnostic tools and mitigation strategies.

Related work has further demonstrated that LLMs are more likely to associate male identities with high-status occupations and leadership roles, while associating female identities with caregiving or subordinate roles (Davis et al., 2022). These biases can persist even when overt gendered terms are removed, indicating that stereotypes are deeply embedded in model representations (Han et al., 2021; Shen et al., 2022). Other studies have highlighted that instruction-tuned models may exhibit amplified gender bias (Dubois et al., 2024; Ferrara, 2023; Ouyang et al., 2022), and that such bias extends beyond English, manifesting across multilingual outputs (Gonen et al., 2022; Barikeri et al., 2021).

Our work builds on these findings by evaluating how gender bias surfaces in downstream political tasks. Unlike prior studies that focus on occupational associations or sentence completions, we examine whether LLMs disproportionately misclassify female Members of European Parliament

---

(MEPs) during gender prediction, and whether gender assumptions influence voting simulation accuracy. By grounding our analysis in real-world political discourse, we contribute novel insights into how gender bias manifests in high-stakes democratic contexts.

## 2.2 Political Leaning of LLMs

A growing body of research has documented that LLMs exhibit consistent political leanings, particularly toward left-of-center or liberal ideologies. Prior work has employed various political orientation tests, including the Political Compass Test (PCT), Pew Research surveys, and the Political Spectrum Quiz, to measure these biases across models (Potter et al., 2024; Bang et al., 2024; Rozado, 2024b; Feng et al., 2023a; Santurkar et al., 2023; Hartmann et al., 2023; Vijay et al., 2024). These studies largely focus on U.S.-centric contexts and have shown that instruction-tuned LLMs tend to demonstrate stronger left-leaning tendencies than their base models.

For example, Hartmann et al. (2023) and Rozado (2024b) found that LLMs exhibit stronger liberal alignment in response to survey-style questions, even when stripped of politically-charged prompts. Similarly, Vijay et al. (2024) demonstrated that LLMs often subtly favor liberal viewpoints, even when instructed to argue from conservative perspectives. Other work has shown that fine-tuning LLMs on partisan data not only shifts their ideological orientation but also degrades their performance in downstream tasks like misinformation detection (Feng et al., 2023a).

While most of these studies rely on multiple-choice surveys or single-turn prompt evaluations, Potter et al. (2024) and Fisher et al. (2024) explore how political bias manifests in more interactive human-LLM dialogues. These works highlight the persuasive effects of politically-biased LLMs on user beliefs, particularly in the context of the 2024 U.S. Presidential election.

Our study diverges from this existing literature in several important ways. First, we shift the geographic and institutional focus to the European Union, introducing not only a non-English but also a multilingual setting that increases the complexity and diversity of the evaluation. Second, rather than relying on survey-style prompts or ideological questionnaires, we assess political leaning through two task-based evaluations: gender prediction and vote simulation. These tasks more closely mirror potential real-world applications of LLMs in political analysis and decision-support systems. Finally, we examine the impact of instruction tuning with political group identifiers on fairness across ideological lines, providing insights into potential mitigation strategies for political bias.

## 3 Data Collection

Debates in the European Parliament take place during plenary sessions, where MEPs deliberate on legislative proposals, reports, and motions. Roll-call votes are a formal voting procedure in which each MEP's vote — 'For', 'Against', or 'Abstain' — is individually recorded and made publicly available. These debates typically precede the vote, offering contextual insights into the positions and arguments put forward by MEPs.

Building on this structure, we introduce EuroParlVote, a novel dataset constructed by collecting roll-call voting records spanning seven years from HowTheyVote.eu (HowTheyVote.eu Team, 2025), covering more than 1,200 MEPs. Using document references provided in the voting metadata, we align these votes with the corresponding debates (Koehn, 2005; Rabinovich et al., 2017; Vanmassenhove and Hardmeier, 2018; Yang et al., 2024, 2023). To ensure relevance of the data, we retain only those debate speeches delivered by MEPs who were present and cast a vote on the associated motion.

We enrich each MEP entry with demographic attributes sourced from their respective Wikipedia pages, including political group affiliation, country, and date of birth, as well as publicly-listed social media accounts (Facebook and Twitter). For gender annotation, we follow a heuristic approach inspired by prior work (Wagner et al., 2016; Reagle and Rhue, 2011): if the English Wikipedia text contains male pronouns (e.g., *he/him/his*), the MEP is labeled as male; if it contains female pronouns (e.g., *she/her*), the label is female. In cases lacking explicit gender indicators, we manually annotate gender based on the MEPs' list pages.[4]

We exclude instances where the vote was marked as 'Abstain', or where either the debate topic or speech was missing. The resulting dataset contains approximately 22K debate speeches linked to 956 unique topics, each paired with MEP-level votes.

We partition the dataset into training, develop-

---

[4] https://www.europarl.europa.eu/meps/en/full-list/all

| Code | Full Name | Political Leaning | Train | Dev | Test |
|------|-----------|-------------------|-------|-----|------|
| GUE/NGL | The Left group in the EP — Nordic Green Left | Far-Left | 1,155 | 133 | 138 |
| GREEN_EFA | Group of the Greens/European Free Alliance | Left | 2,089 | 145 | 140 |
| SD | Group of the Progressive Alliance of Socialists and Democrats | Center-Left | 4,414 | 235 | 207 |
| RENEW | Renew Europe Group | Center / Liberal | 2,826 | 139 | 154 |
| EPP | Group of the European People's Party (Christian Democrats) | Center-Right | 5,004 | 294 | 294 |
| ECR | European Conservatives and Reformists Group | Right | 1,582 | 222 | 250 |
| ID | Identity and Democracy Group | Far-Right | 1,064 | 280 | 248 |
| NI | Unaffiliated Members | Mixed / Variable | 872 | 100 | 117 |

Table 1: Political group codes, full names, ideological positions, and their counts across train/dev/test splits.

| Split | FOR | AGAINST | Male | Female |
|-------|-----|---------|------|--------|
| Train | 16,713 | 2,293 | 55.1% | 44.9% |
| Dev | 774 | 774 | 59.0% | 41.0% |
| Test | 774 | 774 | 59.0% | 41.0% |

Table 2: Vote label counts and gender proportions in the training, development, and test splits.

ment, and test sets using an approximately 8:1:1 ratio. To preserve real-world distributional characteristics, the training set retains the original class imbalance. In contrast, both the development and test sets are balanced across vote labels to support fair evaluation of prediction performance.

Table 1 provides a breakdown of political group affiliations and their ideological positions across splits. Political leanings are determined based on established expert-coded classifications from Parl-Gov and CHES datasets (Döring and Manow, 2023; Polk et al., 2017; Bakker et al., 2015), and validated through European Parliament political group analyses (Hix et al., 2016). Table 2 summarizes the split and gender distribution, it demonstrates a nearly equal gender split across the three sets.

## 4 Investigating Gender Bias of LLMs

### 4.1 Gender Classification Task

The first task involves predicting the gender of the MEPs based on their debate speeches. To accomplish this, we employed the following prompt:

```
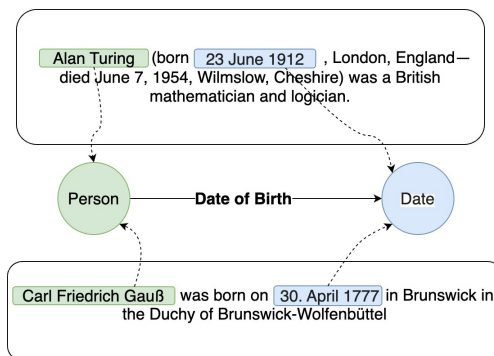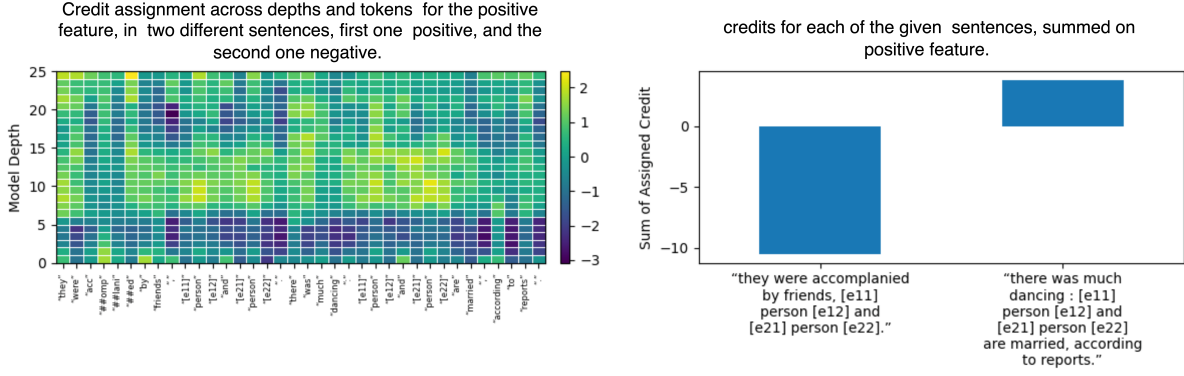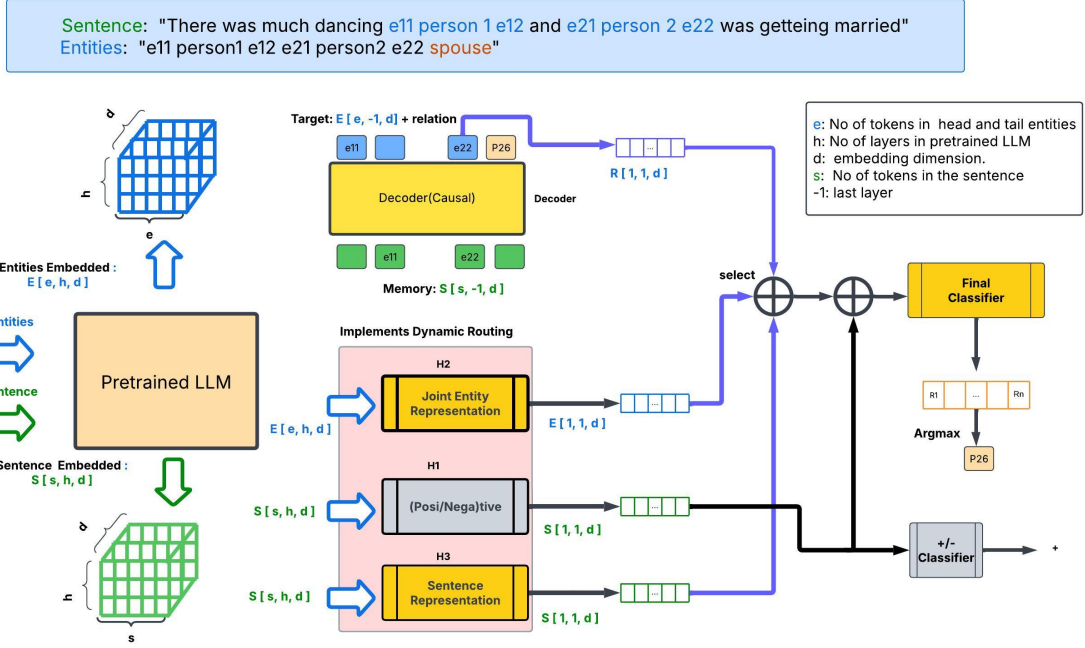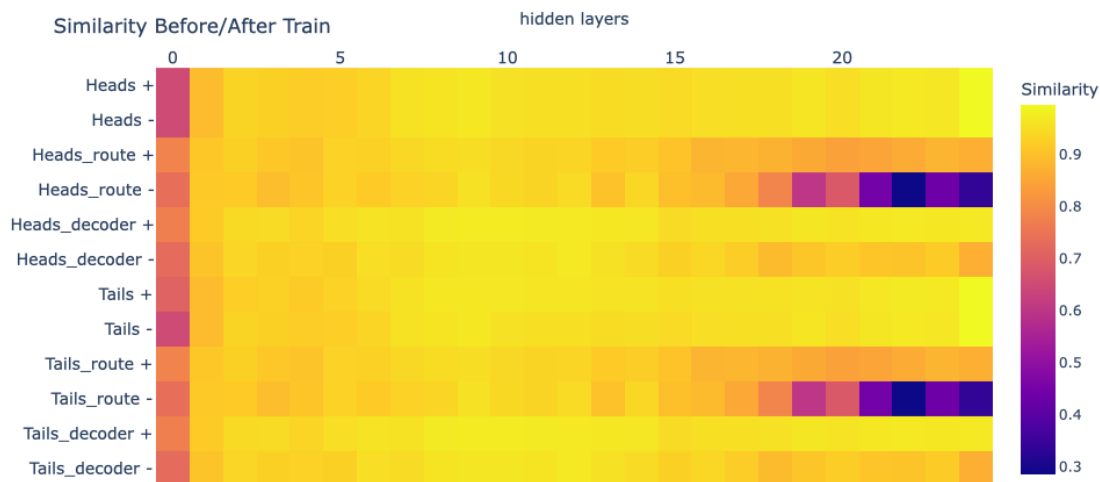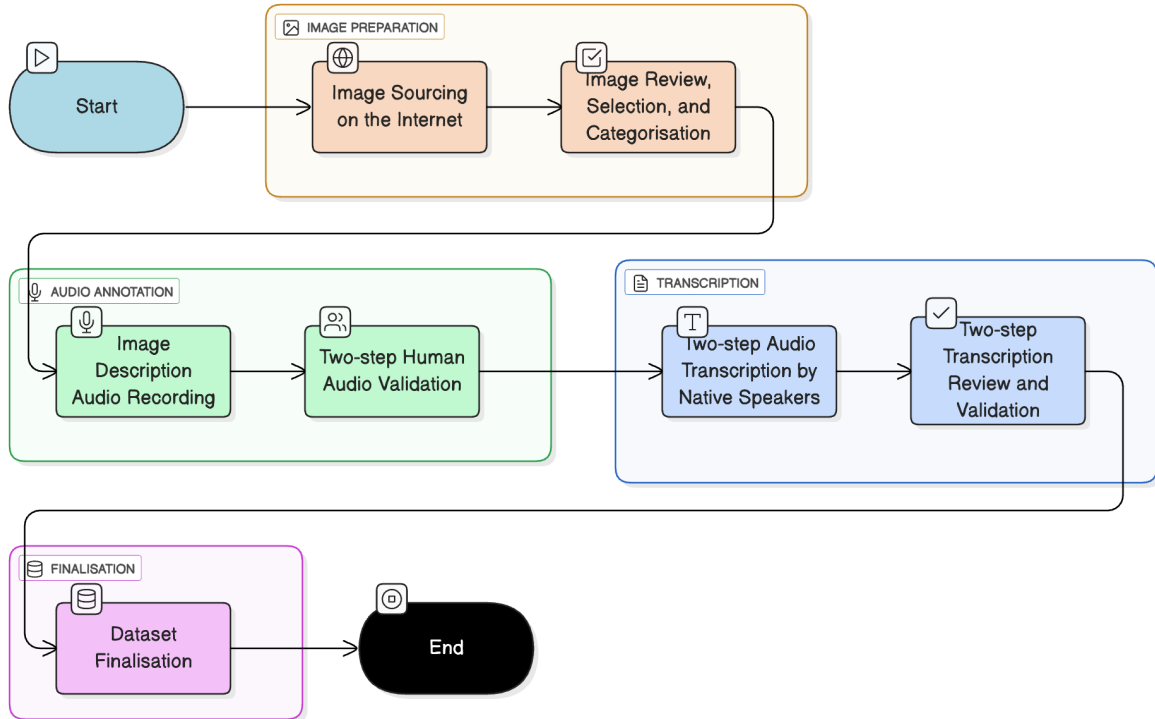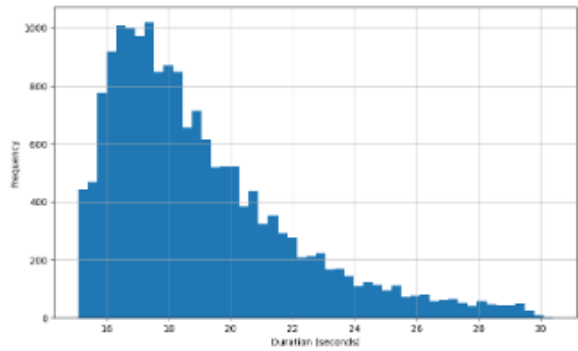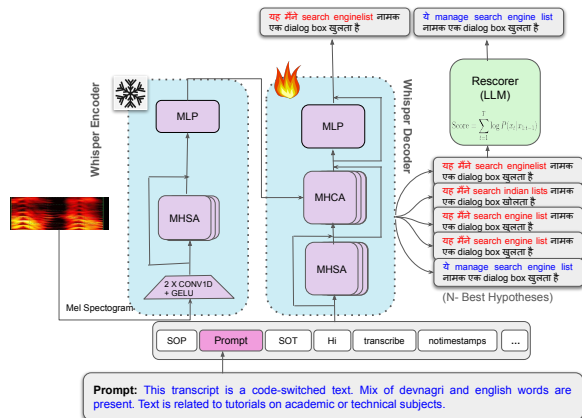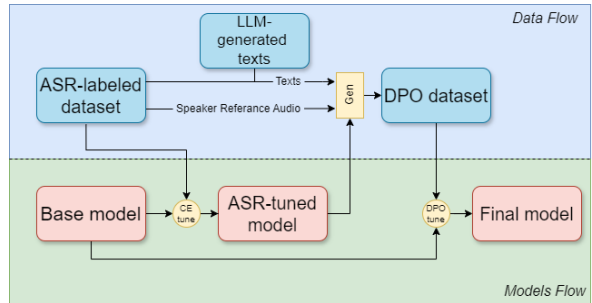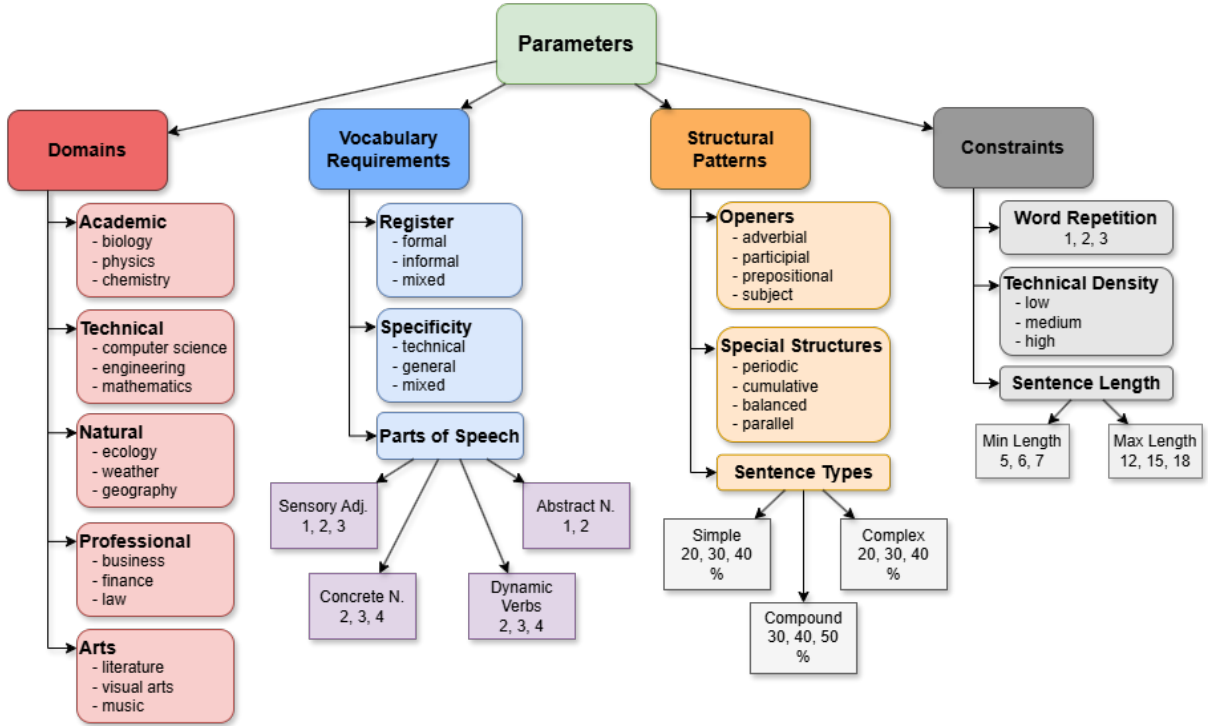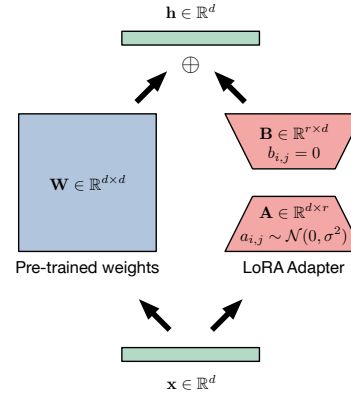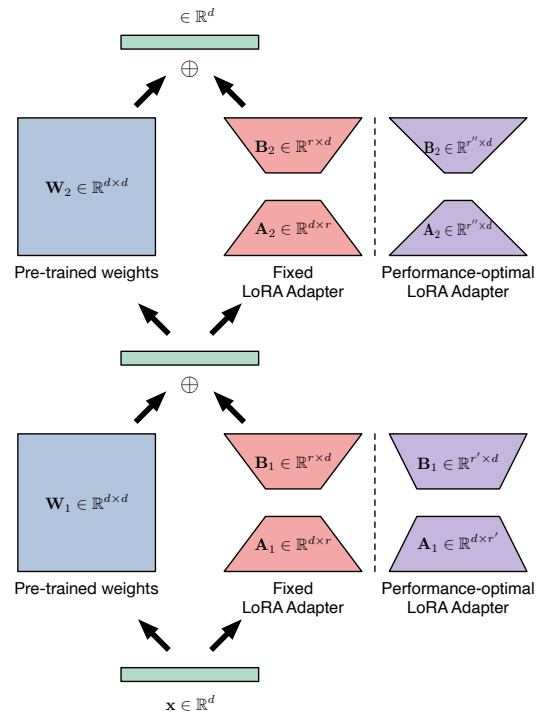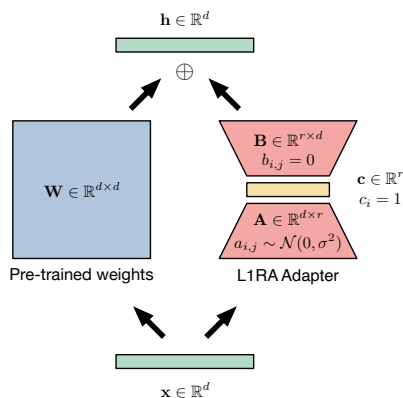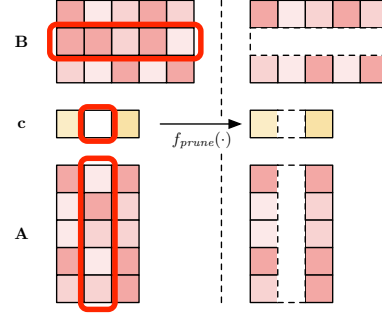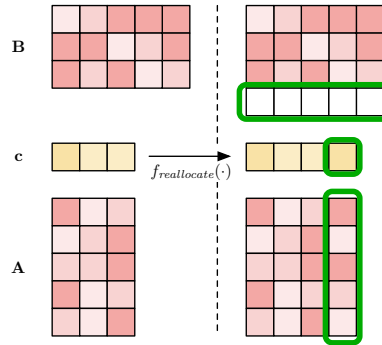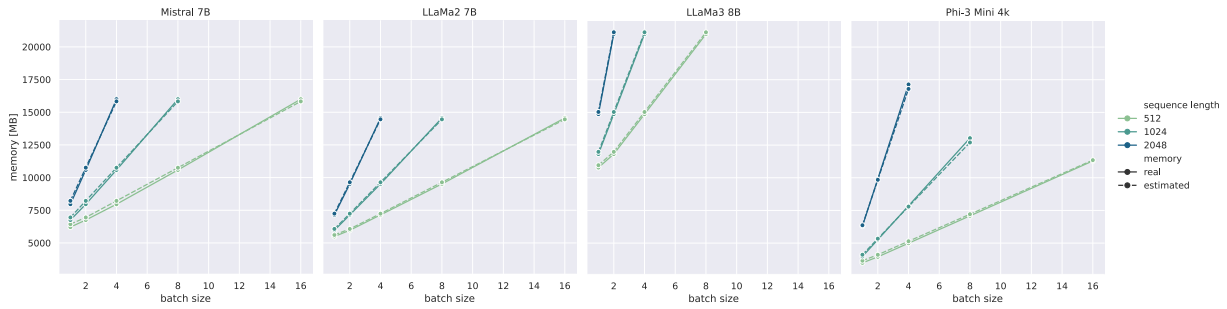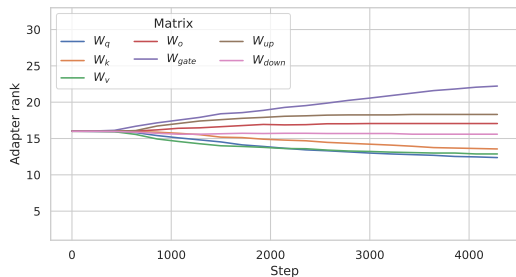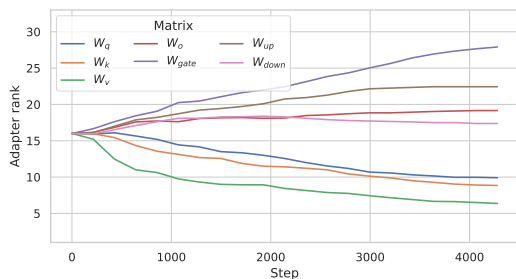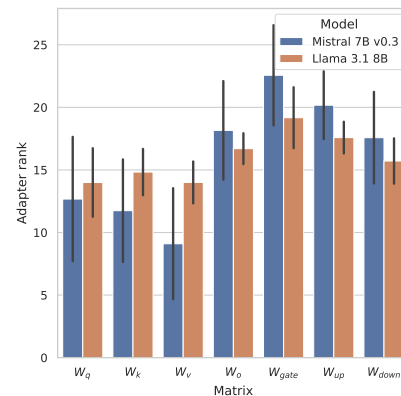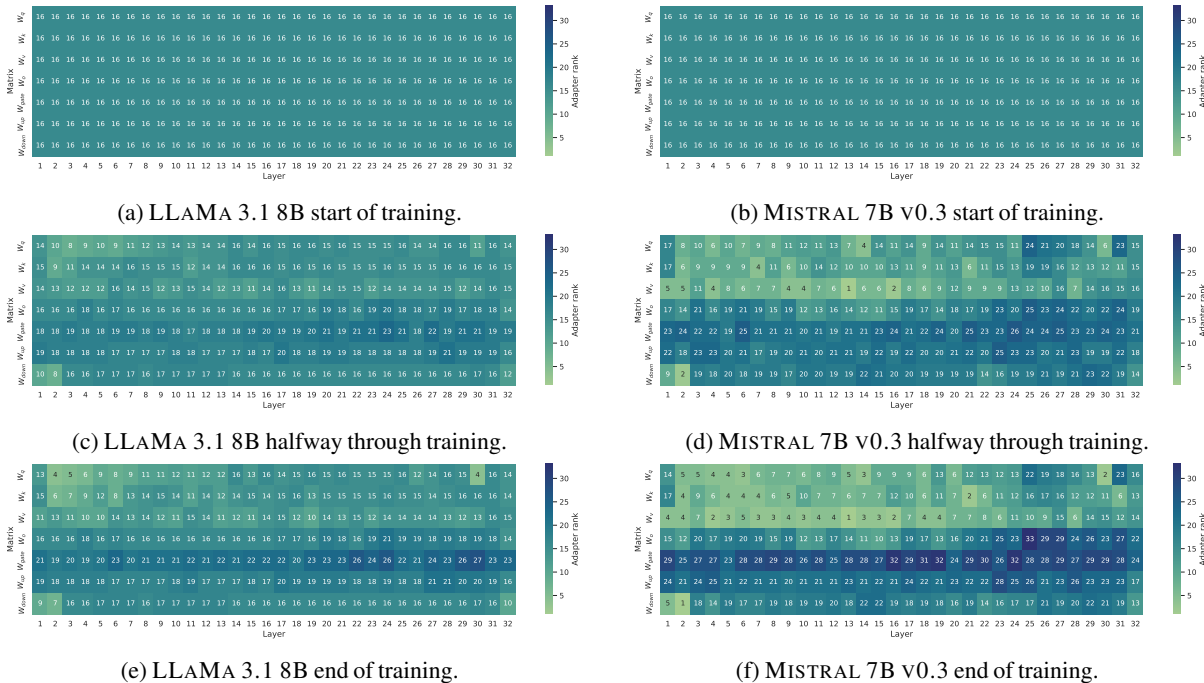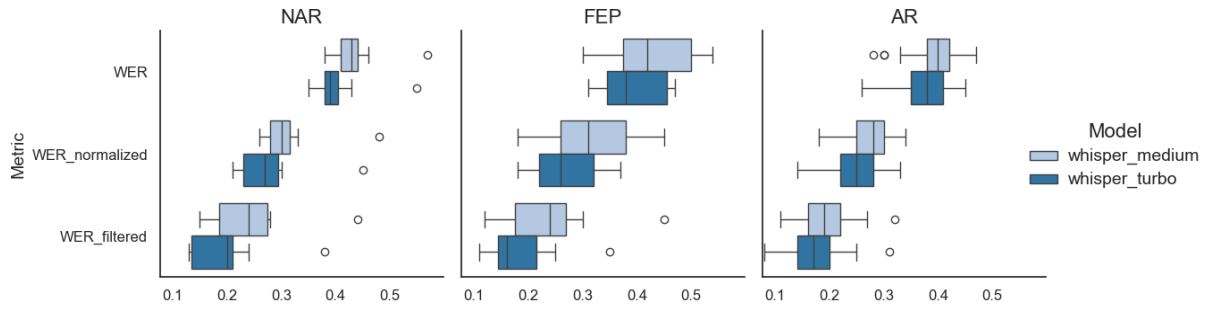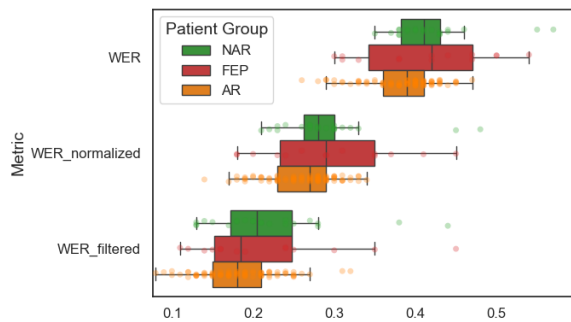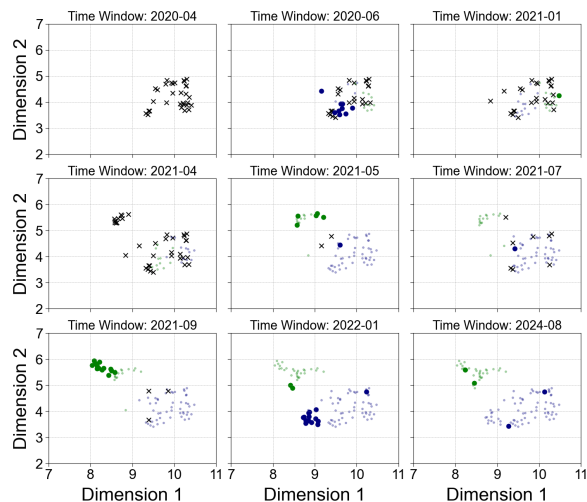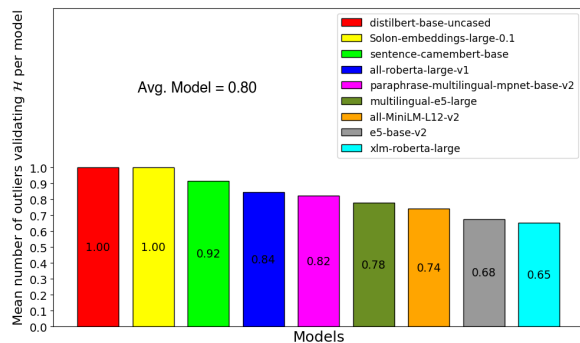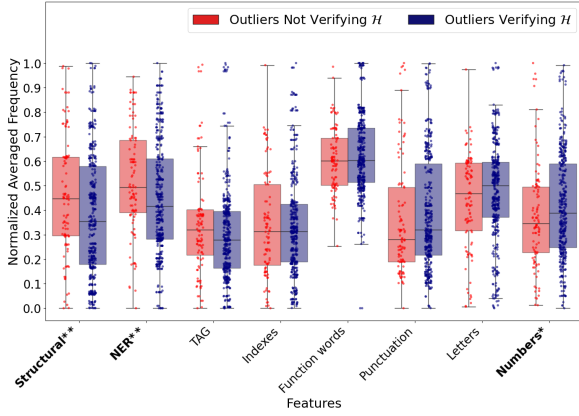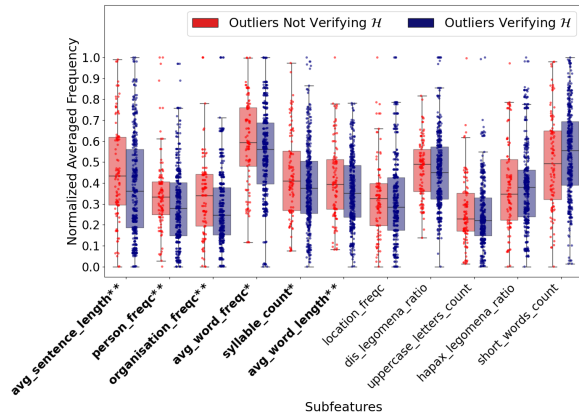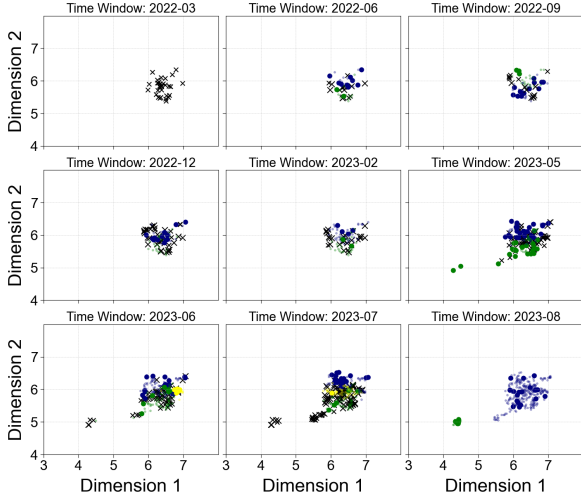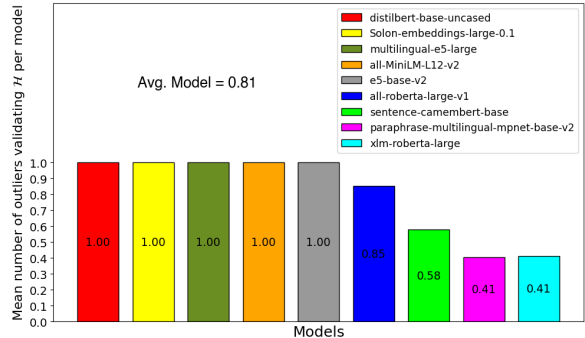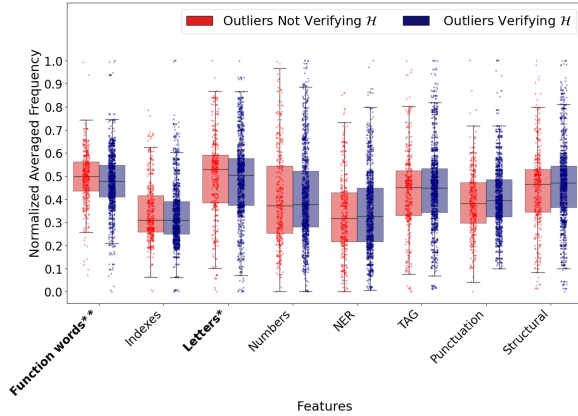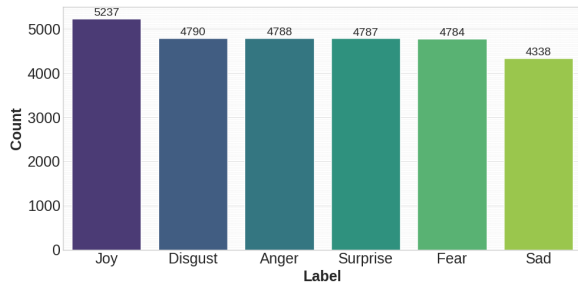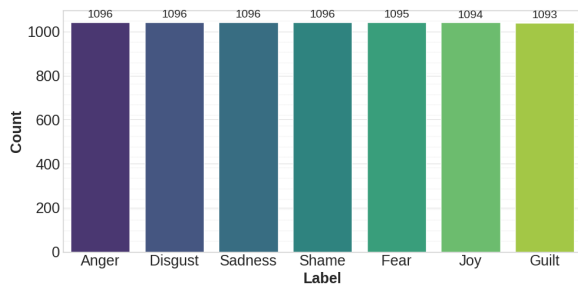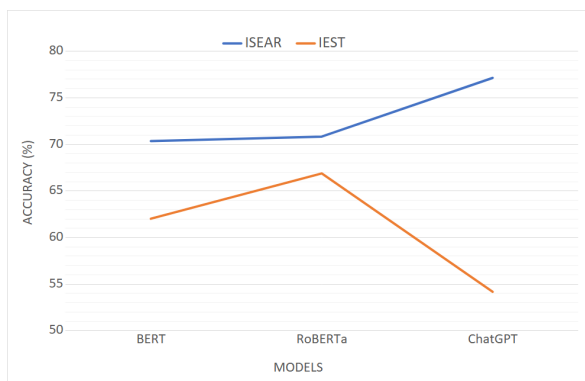Analyze the European Parliament debate speech to determine
whether the speaker is male or female. Please provide: 1.
A gender prediction: "Male" or "Female"; 2. A confidence
score on a scale of 1–5; 3. A rationale for the prediction.
```

We experiment with several language models, including proprietary systems such as GPT-4o (gpt-4o-2024-11-20) (OpenAI, 2024), Gemini-2.5-Flash (gemini-2.5-flash-preview-04-17) (Team et al., 2025; Google, 2025), and Claude-3.5 (claude-

3-5-haiku-20241022) (Anthropic, 2024), as well as open-weight models such as LLaMA-3.2 (LLaMA-3.2-3B-Instruct) (Grattafiori et al., 2024; Touvron et al., 2024) and Mistral-large (mistral-large-2411) (MistralAI, 2023).[5]

Evaluation was conducted in a zero-shot setting on the EuroParlVote test set. Performance on the gender prediction task is reported in terms of Accuracy, F1 scores, and AUC-ROC in Table 3. The AUC-ROC (Area Under the Receiver Operating Characteristic Curve) (Hanley and McNeil, 1982) measures the model's ability to distinguish between male and female classes based on prediction confidence, with higher values indicating better discriminatory capacity.

We summarize our main findings as follows:

**Gender prediction in political discourse is notably more difficult than in other textual domains.** Accuracy across all LLMs ranges from roughly 60 to 70, which is significantly lower than the 80+ accuracy typically observed on the same task in domains such as blogs or news articles (HaCohen-Kerner, 2022; Mukherjee and Liu, 2010). This gap reflects the challenging nature of the EuroParlVote dataset, which features high-register language and rhetorical complexity. As documented in prior sociolinguistic research (Eckert and McConnell-Ginet, 2013; Wodak and Benke, 1997), explicit gender markers are often absent in formal political speech, and the frequent use of irony, sarcasm, and indirect criticism makes gender inference especially difficult.

**Proprietary models outperform open-weight models in both accuracy and calibration.** As shown in Table 3, open-weight models such as LLaMA-3.2 and Mistral-large achieve relatively low accuracy and AUC-ROC scores—LLaMA-3.2,

---

[5]Model details are provided in Appendix Table 10.

| Model | ACC | F1-F | F1-M | AUC |
|---|---|---|---|---|
| LLaMA-3.2-3B | 60.01 | 37.16 | 70.68 | 55.26 |
| Mixtral-Large | 63.60 | 44.85 | 72.84 | 58.50 |
| Claude-3.5 | 64.25 | 60.86 | 67.10 | 63.20 |
| Gemini-2.5-Flash | 65.23 | 51.63 | 72.11 | 55.89 |
| GPT-4o | 61.02 | 59.93 | 62.05 | 66.19 |

Table 3: Gender prediction performance (Accuracy, F1 for Female/Male, AUC-ROC) across models. Lower F1-F scores highlight gender bias toward female MEPs.

for instance, yields 60.01 accuracy and an AUC-ROC of 55.26. In contrast, proprietary models like Claude-3.5, Gemini-2.5 and GPT-4o perform substantially better, indicating its superior capacity to separate classes under uncertainty.

**Proprietary models are also more fair, particularly in their treatment of female speakers.** A clear disparity is observed in the F1 scores for female MEPs, where open-weight models demonstrate a strong male prediction bias. For example, LLaMA-3.2 achieves only 37.16. in F1-Female. As shown in the confusion matrix in Figure 4, surprisingly 71.13% of female speakers are labeled incorrectly as male, compared to just 18.38% misclassification for male speakers. In contrast, GPT-4o shows substantially higher F1-Female scores, indicating a more balanced performance and reduced gender bias.

## 4.2 Voting Simulation Task

The next step in our gender bias investigation involved evaluating how LLMs predict MEP voting positions when presented with debate topics and speeches. Specifically, we tasked the LLMs with simulating whether an MEP would vote with the following prompt:

```
Simulate as a European Parliament MEP. Analyze the debate
topic and speech then state your voting position. Please
provide: 1. Your position ("For" for positive support, or
"Against" for negative rejection) 2. Confidence level on a
scale of 1-5 3. Reasoning for your prediction
```

To explicitly evaluate the impact of demographic cues, we optionally appended the hint: You are a (male|female) MEP at the end of the prompt. This allows the model to consider gender as an explicit feature when making predictions.

We conducted experiments across five settings on LLMs to assess the impact of gender cues: (1) *Without Gender* — no gender information provided; (2) *With Gender* — the prompt included the speaker's actual gender; (3) *All Male* — all

MEPs were labeled as male, regardless of their true gender.; (4) *All Female* — all MEPs were labeled as female; and (5) *Swapped Gender* — Each MEP's gender label was reversed (male ↔ female). This setup enables analysis of how gender signals influence model predictions.

Table 4 summarizes the results across five LLMs. Overall, we observe that injecting gender information leads to nuanced, model-specific effects on performance. For most LLMs, the *All Female* setting resulted in the lowest performance across all metrics—especially in predicting *Against* votes, where F1 scores were notably reduced. Conversely, the *All Male* setting often produced the highest or near-highest results, indicating that male contexts are more aligned with model expectations or learned priors.

Consistent with our findings in Section 4.1, proprietary models not only achieved stronger overall performance but also demonstrated greater fairness. For instance, GPT-4o exhibited minimal variation across gender settings, maintaining relatively high accuracy and AUC-ROC regardless of gender manipulation. This indicates superior robustness to demographic perturbations compared to open-weight and earlier models, which showed more pronounced gender-related performance disparities.

## 4.3 Baselines

To contextualize the learnability of the vote prediction task, we define three baselines representing lower, intermediate, and upper bounds of achievable performance:

**Random baseline** assigns votes uniformly at random, serving as a minimal lower bound with no use of contextual or structural information.

**Group-majority baseline** predicts each MEP's vote based on the most common vote within their political group in the training set. This reflects group-level voting priors without considering the content of debates, providing a simple metadata-based heuristic.

**Intra-group agreement** predicts each MEP's vote based on the majority decision of their political group for that *specific* vote in the test set. For example, if in a given vote 80% of a group's members voted "For", this baseline predicts "For" for all members of that group for that vote. This approach assumes perfect knowledge of group behavior at test time and therefore acts as a soft upper bound, given the high average within-group agreement (95.29%) observed in our dataset.

| Model | Setting | Accuracy | F1-For | F1-Against | AUC-ROC | Avg Confidence |
|-------|---------|----------|--------|------------|---------|----------------|
| Random | Lower Bound | 50.19 | 50.35 | 50.03 | 50.19 | 2.50 |
| Group-Majority | Baseline | 65.28 | 73.65 | 48.96 | 65.25 | 4.03 |
| Intra-Group Agreement | Upper Bound | 88.28 | 89.44 | 86.84 | 87.98 | 4.76 |
| LLaMA-3.2 | Without Gender | 67.10 | 74.55 | 53.85 | 81.88 | 4.00 |
|  | With Gender | 66.47 | 74.26 | 51.91 | 80.49 | 4.00 |
|  | All Male | 67.64 | 73.04 | 54.78 | 81.10 | 3.99 |
|  | All Female | 63.33 | 72.07 | 46.71 | 79.74 | 3.99 |
|  | Swapped Gender | 65.73 | 73.59 | 51.19 | 80.28 | 3.99 |
| Mistral-Large | Without Gender | 75.42 | 76.80 | 73.15 | 82.30 | 4.02 |
|  | With Gender | 76.95 | 78.10 | 75.12 | 83.25 | 4.05 |
|  | All Male | 77.83 | 78.40 | 76.30 | 83.40 | 4.03 |
|  | All Female | 70.64 | 72.91 | 66.87 | 78.12 | 3.91 |
|  | Swapped Gender | 73.50 | 75.10 | 70.25 | 80.45 | 3.94 |
| Claude-3.5 | Without Gender | 80.61 | 81.17 | 80.03 | 85.50 | 4.52 |
|  | With Gender | 82.03 | 82.31 | 81.74 | 86.62 | 4.56 |
|  | All Male | 81.44 | 81.68 | 81.19 | 86.09 | 4.53 |
|  | All Female | 78.34 | 80.12 | 75.31 | 80.28 | 4.51 |
|  | Swapped Gender | 82.12 | 82.30 | 81.93 | 87.07 | 4.53 |
| Gemini-2.5 | Without Gender | 83.10 | 84.10 | 82.05 | 87.90 | 4.25 |
|  | With Gender | 82.78 | 83.65 | 81.90 | 87.50 | 4.28 |
|  | All Male | 82.34 | 83.40 | 80.95 | 87.00 | 4.26 |
|  | All Female | 81.44 | 81.46 | 81.44 | 86.52 | 4.23 |
|  | Swapped Gender | 82.56 | 83.35 | 81.76 | 87.45 | 4.27 |
| GPT-4o | Without Gender | 84.20 | 85.00 | 83.35 | 88.40 | 3.88 |
|  | With Gender | 83.85 | 84.28 | 83.40 | 88.20 | 3.94 |
|  | All Male | 83.72 | 84.14 | 83.29 | 88.49 | 3.93 |
|  | All Female | 83.66 | 84.10 | 83.19 | 87.95 | 3.95 |
|  | Swapped Gender | 83.79 | 84.19 | 83.37 | 88.32 | 3.94 |

Table 4: Voting prediction performance (%) across gender manipulation settings, showing Accuracy, per-class F1, AUC-ROC, and average model confidence. Highlighted rows denote settings where the model struggled the most.

LLMs outperform the random and group-majority baselines, demonstrating that debate content contains useful predictive signals. However, their performance remains below the intra-group agreement baseline, indicating that while LLMs capture informative linguistic patterns, they do not fully replicate structured group voting dynamics.

# 5 Investigating Political Leaning of LLMs

Given its foundation in the political domain, the EuroParlVote dataset naturally assumes that an MEP's political affiliation plays a significant role in shaping their voting behavior. To examine how LLMs respond to this signal, we extend the voting prediction task with two settings: *Without Group* — using the baseline prompt described in Section 4.2, without any mention of political group; and *With Group* — appending the hint *You are a MEP from XX political group* to the end of the prompt. The key findings are as follows:

**Centrist Groups Are Most Accurately Modeled**
Across nearly all models, predictive accuracy peaks for centrist or liberal groups. groups. As shown

in Table 5, the RENEW group consistently yields top scores — e.g., GPT-4o achieves 88.49 accuracy without group information, and 89.93 with it. Mistral and Gemini also perform strongly on RENEW, suggesting its moderate stance is easier for models to simulate. SD (center-left) also yields robust performance, while center-right EPP typically ranks slightly below SD. This trend supports prior work in U.S.-based political modeling, where LLMs tend to exhibit a mild left-leaning bias (Potter et al., 2024; Rozado, 2024a).

**Group Context Boosts Underrepresented Political Groups** As shown in Table 5, explicitly providing political group identity in the prompt improves model performance across most political groups, especially those at the ideological extremes. For instance, LLaMA-3.2 improves on GUE/NGL (far-left) on ID (far-right). Gemini-2.5 also improves considerably on ID. These improvements are especially prominent for politically underrepresented or extreme groups, where models may otherwise struggle to simulate nuanced voting behavior. The addition of group context acts as a

Figure 1: Accuracy of five LLMs across different political groups. The $x$-axis is sorted by the ideological spectrum of the political groups from far-left to far-right.

compensatory fairness signal, particularly benefiting challenging ideological regions.

**Far-Right Groups Are Simulated More Accurately Than Far-Left** Interestingly, the previously observed claim that LLMs exhibit a mild left-leaning bias does not hold at the ideological extremes. As shown in Figure 1 and Table 5, far-right parties such as ID and ECR consistently achieve higher prediction accuracy than their far-left counterparts GUE/NGL and GREEN_EFA. For example, GPT-4o scores 86.07 on ID and 83.78 on ECR, compared to 75.19 on GUE/NGL and 66.21 on GREEN_EFA. This trend holds across other LLMs. These findings suggest that LLMs simulate right-aligned ideological extremes more confidently than left-aligned ones. A possible explanation is that far-right discourse—often more uniform or rhetorically direct—may be easier for LLMs to model, whereas far-left speeches may exhibit greater lexical diversity or abstract reasoning, making them harder to predict from limited input.

## 6 Discussion

### 6.1 Qualitative Analysis of High-Confidence Gender Misclassifications

To better understand the decision patterns of LLMs in gender classification, we conducted a qualitative analysis of high-confidence errors by GPT-4o. Specifically, we examined 200 cases where the model assigned the incorrect gender label with maximum confidence (confidence level = 4), and analyzed its accompanying rationale.

**Stereotypical Language Cues** The model frequently relied on stereotypical associations between tone and gender. For example, assertive, formal, or analytical language was often interpreted

as male: *The text employs a formal, assertive, and analytical tone... suggests a male speaker.*

Similarly, content emphasizing social or environmental concerns was linked with female identity: *Focus on environmental and social issues... associated with female politicians.*

**Political Group Bias** The model also appeared to entangle political ideology with gender assumptions. For instance, far-left MEPs (GUE/NGL) were more likely to be predicted as female due to themes of equity and justice, while conservative MEPs (e.g., ECR) were predicted as male based on critical or structured argumentation—even when incorrect.

**Age Confounds** Older MEPs (age > 70) were disproportionately misclassified. Formal or traditional speech patterns were often read as male-coded, leading to misclassification of several older female MEPs.

### 6.2 Qualitative Analysis of GPT-4o Voting Misclassifications

We conducted a similar error analysis for voting prediction.

**Over-Reliance on Keywords** The model sometimes defaulted to vote predictions based on topic mentions. If a speech referenced climate policy or human rights—topics often associated with FOR votes—it tended to predict approval, even when the speech criticized the specific legislative proposal.

**Surface Sentiment Over Argumentative Stance** GPT-4o often conflated negative sentiment with opposition. For example, speeches that included strong criticisms of implementation or enforcement were misclassified as AGAINST, despite concluding in support: *The implementation has been disappointing and slow. Nevertheless, we must move forward together.* (Predicted: Against, True: For). This reflects a pattern where the model weighs emotional tone over policy alignment.

**Failure to Detect Sarcasm or Irony** In a few speeches, rhetorical devices or sarcastic phrasing led to misclassification. For example, when a speaker said: *Of course, the Commission never makes mistakes.* (Predicted: For, True: Against) The model interpreted literal sentiment and failed to recognize the ironic critique.

| Model | GUE_NGL | GREEN_EFA | SD | RENEW | EPP | ECR | ID |
|---|---|---|---|---|---|---|---|
| *LLaMA-3.2* | | | | | | | |
| w/o group | 47.40 | 61.40 | **88.94** | 86.33 | 81.29 | 53.60 | 60.00 |
| with group | 49.00 | 63.45 | **87.66** | 85.61 | 80.95 | 50.45 | 63.21 |
| *Claude 3.5* | | | | | | | |
| w/o group | 65.86 | 62.47 | 73.92 | **86.34** | 78.33 | 75.68 | 72.04 |
| with group | 66.51 | 65.73 | 75.26 | **84.83** | 80.40 | 81.17 | 76.88 |
| *Mistral-large* | | | | | | | |
| w/o group | 70.03 | 64.56 | 80.14 | **82.54** | 79.20 | 78.84 | 78.88 |
| with group | 71.00 | 66.50 | 81.50 | **84.00** | 83.50 | 82.00 | 79.32 |
| *Gemini 2.5* | | | | | | | |
| w/o group | 68.42 | 64.83 | 76.17 | **80.58** | 79.59 | 76.13 | 80.36 |
| with group | 68.42 | 71.72 | 75.32 | **85.61** | 82.99 | 86.94 | 85.00 |
| *GPT-4o* | | | | | | | |
| w/o group | 75.19 | 66.21 | 88.09 | **88.49** | 82.99 | 83.78 | 86.07 |
| with group | 77.44 | 67.59 | 87.66 | **89.93** | 82.65 | 86.94 | 88.57 |

Table 5: Voting prediction accuracy (%) across political groups for various models. In each row, the highest score is highlighted in bold. Columns are ordered ideologically (left to right) and color-coded from dark blue (far-left) to dark red (far-right), with gray used for center/liberal groups.

### 6.3 Does LoRA Help Mitigate Gender Bias in LLM-based Gender Classification?

Given the observed gender bias in LLM predictions, we investigate whether commonly used fine-tuning techniques, such as supervised fine-tuning (SFT) and Low-Rank Adaptation (LoRA), can mitigate this bias in gender classification tasks. To explore this, we sampled 5,000 examples from the EuroParl-Vote training set, which exhibited a relatively balanced distribution between male and female MEPs, as shown in Table 2.

We applied LoRA fine-tuning to LLaMA3.2-3B and Mistral-Large models, tuning hyperparameters on the development set. The selected hyperparameters included a lora_dropout of 0.05, lora_alpha of 16, a learning rate of 1e-4, and two training epochs. Evaluation was conducted on the test set following the protocol described in Section 4.1.

Table 6 presents the gender prediction performance of the LoRA-finetuned models. For LLaMA-3.2-3B, LoRA yields a slight improvement in overall accuracy and male F1 score. However, it results in a substantial decline in the female F1 score, suggesting a worsening of gender disparity.

A similar decrease is observed for Mistral-large, this trend is further visualized in the confusion matrices shown in Figure 2, where the left panel corresponds to LLaMA3.2 (LoRA) and the right to Mixtral-Large (LoRA). Both models demonstrate strong performance on male MEPs but struggle sig-

nificantly with female MEPs, reinforcing concerns about gender bias.

These findings align with observations by Ding et al. (2024), who report that LoRA does not consistently reduce or exacerbate disparities across demographic subgroups. Our results suggest that while LoRA may enhance general performance, it may also amplify existing gender imbalances unless explicitly addressed.

| LLM | Accuracy | F1-F | F1-M |
|---|---|---|---|
| LLaMA-3.2-3B | 60.01 | 37.16 | 70.68 |
| LLaMA-3.2-3B (LoRA) | 60.70 | 19.94 | 74.88 |
| Mixtral-large | 63.60 | 44.85 | 72.84 |
| Mixtral-large (LoRA) | 61.80 | 32.14 | 75.28 |

Table 6: Gender prediction performance (%): Accuracy, F1 score for Female, and F1 score for Male using original and LoRA fine-tuned open-weight LLMs.



Figure 2: Confusion matrices of gender prediction using LoRA-finetuned models.

| LLM | Accuracy | F1-F | F1-A |
|---|---|---|---|
| LLaMA-3.2-3B (w/o speech) | 50.12 | 66.61 | 0.87 |
| GPT-4o (w/o speech) | 50.39 | 68.67 | 3.03 |

Table 7: Voting prediction performance (Accuracy, F1-For, F1-Against, all in %) of LLaMA3.2 and simulated GPT-4o when the input excludes the speech.

## 6.4 Investigating the Impact of Speech Context in Voting Simulation

To determine whether LLMs are relying on superficial or trivial cues, we conducted an ablation experiment by masking out the debate speeches in the voting simulation task. Instead, we provided only the debate topic and MEP gender. We evaluated one open-weight LLM (LLaMA3.2-3B) and one proprietary model (GPT-4o) under this setup.

Table 7 shows that the accuracy of both models drops to around 50, close to random guessing. Moreover, both models exhibit a strong prediction bias toward the dominant *For* class, resulting in extremely poor F1 scores for the *Against* class.

When compared to the *With Gender* setting in Table 5, where LLaMA3.2 and GPT-4o achieved 66.47 and 83.85 accuracy respectively, this drop highlights the importance of speech context in LLM-based vote prediction. These results confirm that LLMs do not simply rely on gender or group priors but benefit substantially from the semantic content of the debate speeches.

This finding also suggests that using debate speech as a primary input provides richer, non-trivial signals for political decision modeling, reinforcing the critical role of context in socially grounded LLM applications.

## 6.5 Investigating the Limitation of Machine Translation on Voting Prediction

Given the multilingual nature of the EuroParlVote dataset, all results reported in Section 4 and Section 5 have used speeches in their original language. Meanwhile, we were curious whether the originality of the language affects model performance in downstream tasks. This question aligns with concerns raised in prior work on language bias in multilingual NLP systems (Yang et al., 2024).

To investigate this, we translated the speeches in the test set using three methods: GPT-4o, T5(Raffel et al., 2020), and the Google Translate API(Google, 2024). We then used the best-performing model, GPT-4o, to replicate the voting prediction experiment described in Section 4.2, under the setting without gender or group metadata.

As shown in Table 8, all translated versions yield lower accuracy than the original-language speeches. This result is consistent with expectations, as translation may introduce noise or omit important contextual signals. It also highlights the value and authenticity of our benchmark, which retains original native-language inputs.

| Translator | Accuracy | F1-F | F1-A |
|---|---|---|---|
| GPT-4o | 78.10 | 80.12 | 75.44 |
| T5 | 75.84 | 78.65 | 72.03 |
| Google API | 76.35 | 79.02 | 72.88 |
| No translation | **84.20** | **85.00** | **83.35** |

Table 8: Voting prediction performance (Accuracy, F1-For, F1-Against, all in %) using translated speeches and original speeches with GPT-4o as the prediction model.

## 7 Conclusion

We introduced EuroParlVote, a benchmark dataset aligning MEP debate speeches with roll-call votes and demographic metadata, enabling fine-grained evaluation of LLMs across gender and political group dimensions in a real-world democratic setting.

Our findings reveal persistent *gender* and *ideological* biases in current LLMs. Proprietary models such as GPT-4o show greater robustness and fairness, while open-weight models like LLaMA-3.2 benefit from explicit contextual cues (e.g., political group identifiers) but still fail in predictable ways, including over-reliance on sentiment polarity, misinterpreting hedging or irony, and insufficiently integrating context or speaker intent. In the futuer work, incorporating discourse signals such as group alignment, prior voting records, and procedural vs. policy distinctions may improve robustness.

We also find that LoRA fine-tuning fails to mitigate gender disparities, and that translating multilingual debates reduces performance—underscoring the importance of native-language inputs. To our knowledge, this is the first study to jointly examine gender and political fairness in LLMs in a multilingual parliamentary context. We hope EuroParlVote fosters research into the socio-political implications of LLM deployment and encourages the development of fairer, more context-aware NLP systems for political applications.

## Limitations

This study has several limitations. First, due to budget constraints and limited API access, we did not conduct ablation studies across all trending model variants or include very large-scale LLMs (e.g., LLaMA-3.3-70B). Instead, we selected a diverse yet manageable set of proprietary and open-weight models to facilitate consistent, cross-comparative analysis. Second, our focus was on identifying and analyzing bias rather than developing or fine-tuning mitigation techniques, which typically require additional training cycles, labeled data, or access to model internals—challenges that are particularly acute for proprietary models. Lastly, both the content of European Parliament debates and the capabilities of LLMs are dynamic and evolving. As a result, our findings may not fully generalize to future model versions or accurately reflect shifts in political discourse and societal context.

## Acknowledgments

## References

Anthropic. 2024. Claude 3.5 haiku model card. `https://www.anthropic.com/claude/haiku`. Introduces Claude 3.5Haiku as the fastest model in the Claude 3 family, launched October 22, 2024.

Ryan Bakker, Catherine de Vries, Erica Edwards, Liesbet Hooghe, Seth Jolly, Gary Marks, Jonathan Polk, Jan Rovny, Marco Steenbergen, and Milada Vachudova. 2015. Measuring party positions in europe: The chapel hill expert survey trend file, 1999–2010. *Party Politics*, 21(1):143–152.

Yejin Bang, Delong Chen, Nayeon Lee, and Pascale Fung. 2024. Measuring political bias in large language models: What is said and how it is said. In *Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 11142–11159, Bangkok, Thailand. Association for Computational Linguistics.

Shikha Barikeri, Vinay Uday Prabhu, and 1 others. 2021. Redditbias: A real-world resource for bias evaluation and debiasing of language models. In *Proceedings of the 2021 Conference on Fairness, Accountability, and Transparency (FAccT)*.

Paul Christiano, Jan Leike, Tom B Brown, Miljan Martic, Shane Legg, and Dario Amodei. 2017. Deep reinforcement learning from human preferences. *arXiv preprint arXiv:1706.03741*.

Sara R. Davis, Cody J. Worsnop, and Emily M. Hand. 2022. Gender bias recognition in political news articles. *Machine Learning with Applications*, 8:100304.

Zhoujie Ding, Yifan Wang, Yuchen Zhang, Yao Li, and William Yang Wang. 2024. On the fairness of low-rank adaptation for large language models. In *Proceedings of the Conference on Language Modeling (COLM)*.

Yann Dubois, Xuechen Li, Rohan Taori, Tianyi Zhang, Ishaan Gulrajani, Jimmy Ba, Carlos Guestrin, Percy Liang, and Tatsunori B. Hashimoto. 2024. Alpaca-farm: A simulation framework for methods that learn from human feedback. *Preprint*, arXiv:2305.14387.

Holger Döring and Philip Manow. 2023. Parlgov: A database of political parties, elections and governments. Accessed: 2024-07-14.

Penelope Eckert and Sally McConnell-Ginet. 2013. *Language and Gender*. Cambridge University Press.

Shangbin Feng, Chan Young Park, Yuhan Liu, and Yulia Tsvetkov. 2023a. From pretraining data to language models to downstream tasks: Tracking the trails of political biases leading to unfair NLP models. In *Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (ACL 2023)*, pages 11737–11762, Toronto, Canada. Association for Computational Linguistics.

Shi Feng, Heejin Kim, Daniel DeBlasio, Aman Madaan, Graham Neubig, and Ximing Liu. 2023b. How does pretraining data affect alignment in large language models? In *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing (EMNLP)*.

Emilio Ferrara. 2023. Should chatgpt be biased? challenges and risks of bias in large language models. *First Monday*.

Jillian Fisher, Shangbin Feng, Robert Aron, Thomas Richardson, Yejin Choi, Daniel W. Fisher, Jennifer Pan, Yulia Tsvetkov, and Katharina Reinecke. 2024. Biased ai can influence political decision-making. *arXiv preprint arXiv:2410.06415*.

Hila Gonen, Shauli Ravfogel, and Yoav Goldberg. 2022. Analyzing gender representation in multilingual models. In *Proceedings of the 7th Workshop on Representation Learning for NLP*, pages 67–77, Dublin, Ireland. Association for Computational Linguistics.

Google. 2024. Google cloud translation api. `https://cloud.google.com/translate`. Accessed: 2024-05-17.

Google. 2025. Start building with gemini 2.5 flash. Google Developers Blog. Preview release via Gemini API in Google AI Studio and Vertex AI.

Aaron Grattafiori, Abhimanyu Dubey, Abhinav Jauhri, Abhinav Pandey, Abhishek Kadian, Ahmad Al-Dahle, Aiesha Letman, Akhil Mathur, Alan Schelten, Alex Vaughan, Amy Yang, Angela Fan, Anirudh Goyal, Anthony Hartshorn, Aobo Yang, Archi Mitra, Archie Sravankumar, Artem Korenev, Arthur Hinsvark, and 542 others. 2024. The llama 3 herd of models. *Preprint*, arXiv:2407.21783.

Yaakov HaCohen-Kerner. 2022. Survey on profiling age and gender of text authors. *Expert Systems with Applications*, 199:117140.

Xudong Han, Timothy Baldwin, and Trevor Cohn. 2021. Diverse adversaries for mitigating bias in training. In *Proceedings of the 16th Conference of the European Chapter of the Association for Computational Linguistics: Main Volume*, pages 2760–2765.

James A Hanley and Barbara J McNeil. 1982. The meaning and use of the area under a receiver operating characteristic (roc) curve. *Radiology*, 143(1):29–36.

Jochen Hartmann, Jasper Schwenzow, and Maximilian Witte. 2023. The political ideology of conversational AI: Converging evidence on chatgpt's pro-environmental, left-libertarian orientation. *arXiv preprint arXiv:2301.01768*.

Simon Hix, Abdul Noury, and Gérard Roland. 2016. *Democratic Politics in the European Parliament*. Cambridge University Press, Cambridge.

HowTheyVote.eu Team. 2025. Howtheyvote.eu: European parliament roll-call vote transparency. Accessed: 2025-05-14.

Edward J. Hu, Yelong Shen, Phillip Wallis, Zeyuan Allen-Zhu, Yuanzhi Li, Shean Wang, Lu Wang, and Weizhu Chen. 2021. Lora: Low-rank adaptation of large language models. *arXiv preprint arXiv:2106.09685*.

Philipp Koehn. 2005. Europarl: A parallel corpus for statistical machine translation. In *Proceedings of the 10th Machine Translation Summit*, pages 79–86, Phuket, Thailand.

Hadas Kotek, Rikker Dockum, and David Q. Sun. 2023. Gender bias and stereotypes in large language models. In *Proceedings of the ACM Collective Intelligence Conference (CI '23)*, pages 12–24. ACM.

MistralAI. 2023. Mixtral of experts. `https://mistral.ai/news/mixtral-of-experts/`. Accessed: 2024-05-14.

Arjun Mukherjee and Bing Liu. 2010. Improving gender classification of blog authors. In *Proceedings of the 2010 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pages 207–217, Cambridge, MA, USA. Association for Computational Linguistics.

OpenAI. 2024. Gpt-4o system card. `https://openai.com/index/gpt-4o-system-card/`. System card documenting capabilities, limitations, and safety evaluations of GPT-4o.

Long Ouyang, Jeffrey Wu, Xu Jiang, Diogo Almeida, and 1 others. 2022. Training language models to follow instructions with human feedback. *arXiv preprint arXiv:2203.02155*.

Jonathan Polk, Jan Rovny, Ryan Bakker, Erica Edwards, Liesbet Hooghe, Seth Jolly, Jelle Koedam, Filip Kostelka, Gary Marks, Gijs Schumacher, Marco Steenbergen, Milada Vachudova, and Marko Zilovic. 2017. Explaining the salience of anti-elitism and reducing political corruption for political parties in europe with the 2014 chapel hill expert survey data. *Research & Politics*, 4(1):1–9.

Yujin Potter, Shiyang Lai, Junsol Kim, James Evans, and Dawn Song. 2024. Hidden persuaders: Llms' political leaning and their influence on voters. *arXiv preprint arXiv:2410.24190*. Presented at EMNLP 2024.

Ella Rabinovich, Raj Nath Patel, Shachar Mirkin, Lucia Specia, and Shuly Wintner. 2017. Personalized machine translation: Preserving original author traits. In *Proceedings of the 15th Conference of the European Chapter of the Association for Computational Linguistics: Volume 1, Long Papers*, pages 1074–1084, Valencia, Spain. Association for Computational Linguistics.

Colin Raffel, Noam Shazeer, Adam Roberts, Katherine Lee, Sharan Narang, Michael Matena, Yanqi Zhou, Wei Li, and Peter J. Liu. 2020. Exploring the limits of transfer learning with a unified text-to-text transformer. *Journal of Machine Learning Research*, 21(140):1–67.

Joseph Reagle and Lauren Rhue. 2011. Gender bias in wikipedia and britannica. *International Journal of Communication*, 5:1138–1158.

David Rozado. 2024a. Political bias in language models: Evidence from gpt-4 and claude. *SSRN*. Available at `https://ssrn.com/abstract=4670854`.

David Rozado. 2024b. The political preferences of llms. *arXiv preprint arXiv:2402.01789*.

Shibani Santurkar, Esin Durmus, Faisal Ladhak, Cinoo Lee, Percy Liang, and Tatsunori Hashimoto. 2023. Whose opinions do language models reflect? In *Proceedings of the 40th International Conference on Machine Learning*, volume 202 of *Proceedings of Machine Learning Research*, pages 29971–30004. PMLR.

Shibani Santurkar, Toniann Pitassi, Stefano Ermon, and Deep Ganguli. 2024. Whose opinions do language models reflect? In *Proceedings of the 2024 International Conference on Learning Representations (ICLR)*.

Aili Shen, Xudong Han, Trevor Cohn, Timothy Baldwin, and Lea Frermann. 2022. Optimising equal opportunity fairness in model training. In *Proceedings of the 2022 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pages 4073–4084.

Gemini Team, Rohan Anil, Sebastian Borgeaud, Jean-Baptiste Alayrac, Jiahui Yu, Radu Soricut, Johan Schalkwyk, Andrew M. Dai, Anja Hauth, Katie Millican, David Silver, Melvin Johnson, Ioannis Antonoglou, Julian Schrittwieser, Amelia Glaese, Jilin Chen, Emily Pitler, Timothy Lillicrap, Angeliki Lazaridou, and 1332 others. 2025. Gemini: A family of highly capable multimodal models. *Preprint*, arXiv:2312.11805.

Hugo Touvron, Thibaut Lavril, Gautier Izacard, Xavier Martinet, and 1 others. 2024. Llama 3: Open foundation and instruction-tuned models. https://ai.meta.com/blog/meta-llama-3/. Meta AI.

Eva Vanmassenhove and Christian Hardmeier. 2018. Europarl datasets with demographic speaker information. In *Proceedings of the 21st Annual Conference of the European Association for Machine Translation*, Alacant, Spain. European Association for Machine Translation.

Supriti Vijay, Aman Priyanshu, and Ashique R. KhudaBukhsh. 2024. When neutral summaries are not that neutral: Quantifying political neutrality in llm-generated news summaries. *arXiv preprint arXiv:2410.09978*.

Claudia Wagner, Eduardo Graells-Garrido, David Garcia, and Filippo Menczer. 2016. Women through the glass ceiling: Gender asymmetries in wikipedia. *EPJ Data Science*, 5(1):5.

Stefanie Walter, Lucy Kinski, and Zsófia Boda. 2023. Who talks to whom? using social network models to understand debate networks in the european parliament. *European Union Politics*, 24(2):410–423.

Ruth Wodak and G. Benke. 1997. *Gender as a sociolinguistic variable: New perspectives on variation studies.*, pages 127–150. Blackwell.

Jinrui Yang, Timothy Baldwin, and Trevor Cohn. 2023. Multi-EuP: The multilingual European parliament dataset for analysis of bias in information retrieval. In *Proceedings of the 3rd Workshop on Multi-lingual Representation Learning (MRL)*, pages 282–291, Singapore. Association for Computational Linguistics.

Jinrui Yang, Fan Jiang, and Timothy Baldwin. 2024. Language bias in multilingual information retrieval: The nature of the beast and mitigation methods. In *Proceedings of the Fourth Workshop on Multilingual Representation Learning (MRL 2024)*, pages 280–292, Miami, Florida, USA. Association for Computational Linguistics.

# A Appendix

This appendix provides additional details about the EuroParlVote dataset. Figure 3 visualizes the demographic distributions across the training, development, and test splits, highlighting attributes such as gender, political group, age, country, and vote label. Table 9 presents the full mapping of country codes, ISO Alpha-2 codes, and the number of examples per country across the three dataset splits. Notably, while the United Kingdom (GBR) is no longer an EU member, it remains in the dataset due to its historical participation during the data collection period.



Figure 3: Demographic distributions across the training, evaluation, and test sets in the EuroParlVote dataset. Each row shows the distribution by gender, political group (left-to-right political leaning order), age, country, and vote label, respectively.

| Code | ISO Alpha-2 | Country Name | Train | Dev | Test |
|------|-------------|--------------|-------|-----|------|
| AUT | AT | Austria | 562 | 42 | 46 |
| BEL | BE | Belgium | 633 | 49 | 47 |
| BGR | BG | Bulgaria | 417 | 22 | 36 |
| CYP | CY | Cyprus | 204 | 16 | 18 |
| CZE | CZ | Czechia | 531 | 50 | 45 |
| DEU | DE | Germany | 1123 | 97 | 91 |
| DNK | DK | Denmark | 375 | 31 | 32 |
| ESP | ES | Spain | 1044 | 83 | 88 |
| EST | EE | Estonia | 218 | 18 | 21 |
| FIN | FI | Finland | 368 | 28 | 33 |
| FRA | FR | France | 1087 | 89 | 86 |
| GBR | GB | United Kingdom | 946 | 84 | 76 |
| GRC | GR | Greece | 641 | 46 | 47 |
| HRV | HR | Croatia | 323 | 26 | 24 |
| HUN | HU | Hungary | 428 | 32 | 34 |
| IRL | IE | Ireland | 312 | 21 | 23 |
| ITA | IT | Italy | 1057 | 88 | 81 |
| LTU | LT | Lithuania | 232 | 20 | 19 |
| LUX | LU | Luxembourg | 179 | 12 | 15 |
| LVA | LV | Latvia | 200 | 16 | 18 |
| MLT | MT | Malta | 157 | 12 | 13 |
| NLD | NL | Netherlands | 610 | 44 | 49 |
| POL | PL | Poland | 730 | 57 | 60 |
| PRT | PT | Portugal | 504 | 41 | 37 |
| ROU | RO | Romania | 537 | 42 | 44 |
| SVK | SK | Slovakia | 294 | 25 | 21 |
| SVN | SI | Slovenia | 278 | 19 | 21 |
| SWE | SE | Sweden | 493 | 38 | 35 |

Table 9: Mapping of country codes, ISO Alpha-2 codes, country names, and number of examples in Train, Dev, and Test splits in the EuroParlVote dataset. Note that the United Kingdom (GBR) appears due to legacy participation, though it is no longer an EU member.

## B  LLMs Model Configures

This appendix provides additional information on the LLMs evaluated and their performance characteristics. Table 10 summarizes the LLMs used in our experiments, including release dates, parameter sizes, and approximate API pricing.

| LLM | Release Date | Parameters | API Pricing (USD) |
|---|---|---|---|
| LLaMA-3.2-3B-Instruct (Meta) | Sept 2024 | 3.2B | Free for research |
| Mistral-large-2411 (Mistral) | Nov 2024 | 123B | $\sim$1–3/M tokens via Vertex |
| GPT-4o (OpenAI) | Nov 2024 | Not disclosed | 2.50/M input, 5–10/M output |
| Gemini 2.5 Flash (Google) | Apr 2025 | Not disclosed | $\sim$0.26/M tokens (combined est.) |
| Claude 3.5 Haiku (Anthropic) | Oct 2024 | Not disclosed | 0.80/M input, 4.00/M output |

Table 10: Summary of LLMs used in this work, release dates, parameter sizes, and approximate API pricing.

# C  Supplementary Evaluation Details

Figures 4–7 present confusion matrices and confidence distributions for both gender and vote prediction tasks. These visualizations illustrate model behavior across classes and offer insight into confidence calibration and classification asymmetries.



Figure 4: LLMs' Confusion Matrix for gender prediction based on debate speeches.

Figure 5: Distribution of LLMs confidence scores for gender predictions. The violin plot shows the density and spread of confidence levels (1 = low, 5 = high) across predicted genders. Inner boxes indicate the interquartile range and median.



Figure 6: LLMs' Confusion Matrix for vote prediction based on debate topic and debate speeches.

Figure 7: Distribution of LLMs confidence scores for vote predictions grouped by ground-truth label (FOR vs. AGAINST). Wider violins indicate greater density of predictions at that confidence level.

## D   Qualitative Prediction Examples

To further illustrate LLM behavior, we provide qualitative examples of vote and gender prediction outputs for each LLM evaluated. These examples include the ground truth, model predictions, confidence scores, and the corresponding MEP debate speeches. Each instance is selected to highlight notable patterns—such as confident misclassifications, group-specific disparities, or reasoning errors—across different models. The following subsections present examples for LLaMA-3.2, GPT-4o, Gemini 2.5 Flash, Mistral-large, and Claude 3.5 Haiku.

### LLaMA-3.2 Vote Prediction Error Example

---

**Example:** *An EU strategy to reduce methane emissions (debate)*

**Ground Truth:** AGAINST          **LLaMA-3.2 Prediction:** FOR (Confidence: 4)

**MEP Speech**

– Signora Presidente, signora Commissaria, onorevoli colleghi, il metano è un gas serra con un potenziale di riscaldamento globale molto superiore all'anidride carbonica e quindi rappresenta una minaccia importante per l'ambiente. È quindi positivo che la Commissione abbia adottato una strategia per la riduzione delle emissioni di metano. Tuttavia, dobbiamo anche assicurarci che tale strategia non imponga oneri sproporzionati alle piccole e medie imprese agricole, che rischiano di essere danneggiate da regolamenti troppo stringenti.

**LLaMA-3.2 Reasoning**

The MEP speech highlights the importance of the EU's strategy to reduce methane emissions and acknowledges its environmental relevance. The speaker expresses support for the initiative while urging caution to avoid overburdening small farmers. The overall tone is constructive, suggesting alignment with the strategy, leading to a 'FOR' prediction.

---

### LLaMA-3.2 Gender Prediction Error Example

---

**Example:** *2019 Discharge (debate)*

**Ground Truth:** FEMALE          **LLaMA-3.2 Prediction:** MALE (Confidence: 4)

**MEP Speech**

Now a few words on each of these institutions in my role as rapporteur. For the European Commission, I welcomed improvements in budget implementation and the increasing use of performance-based budgeting. However, I also expressed concerns about transparency and accountability in some external aid projects. It is vital that European taxpayers see a return on their investment and that funds are used effectively.

**LLaMA-3.2 Reasoning**

The speaker uses technical language, discusses financial oversight, and presents a structured critique. The formal tone and lack of gender-specific cues may have led the model to predict the speaker as male.

---

**GPT-4o Vote Prediction Error Example**

> **Example:** *The proposed repeal of the law banning female genital mutilation (debate)*
>
> | **Ground Truth:** FOR | **GPT-4o Prediction:** AGAINST (Confidence: 5) |
> | --- | --- |
>
> ### MEP Speech
>
> Wir haben bestehende Konventionen zur Eliminierung von Diskriminierung gegen Frauen – die Istanbul-Konvention, die Genfer Konvention, die UN-Kinderrechtskonvention – alle verurteilen weibliche Genitalverstümmelung eindeutig. Trotzdem erleben wir, dass in bestimmten Ländern über die Abschaffung solcher Gesetze diskutiert wird. Dies ist eine eklatante Verletzung grundlegender Menschenrechte und wir als EU-Parlament müssen klar Stellung beziehen. Es ist unsere Verantwortung, den Schutz von Frauen und Mädchen nicht nur zu fordern, sondern auch aktiv zu verteidigen. Deshalb unterstütze ich diese Resolution mit voller Überzeugung.
>
> ### GPT-4o Reasoning
>
> The MEP speech strongly condemns the practice of female genital mutilation (FGM) and emphasizes existing international conventions against such practices, highlighting them as outdated and not in the interest of women. The speaker calls for the Gambian government to retract its intention to repeal the ban, indicating a clear stance against the repeal. Given the EU's historical position on human rights and gender equality, it is highly likely that the European Parliament would vote against any measure perceived as regressive on these issues. The strong language used in the speech, coupled with the EU's commitment to upholding international human rights standards, suggests a very high probability that the vote will be against the repeal.

**GPT-4o Gender Prediction Error Example**

> **Example:** *Iran, notably the situation of women's rights defenders and imprisoned EU dual nationals*
>
> | **Ground Truth:** MALE | **GPT-4o Prediction:** FEMALE (Confidence: 4) |
> | --- | --- |
>
> ### MEP Speech
>
> – Mr President, I just wanted to concur with what has been said by many speakers. The situation of women in Iran is unacceptable, and particularly from the point of view of the European Union, when I was a member of the Iranian delegation in the last mandate, we worked very hard to keep the nuclear accord in place, and the European Union stuck its neck out in that regard, but that must be reciprocated. As has been said, by particularly allowing women the freedom they deserve, and as one MEP said, men in Iran should not be telling women what to wear. I also concur with my colleague, Clare Daly, on the obligations of the United States in relation to Iran, and the nuclear accord. It was they who pulled out of it, and that was done despite the fact that there was an agreement done with the previous president, President Obama, so that in itself is also a cause for condemnation.

**Gemini 2.5 Vote Prediction Error**

**Example: *Carbon Border Adjustment Mechanism (CBAM)***

**Ground Truth:** AGAINST          **Gemini 2.5 Prediction:** FOR (Confidence: 4)

**MEP Speech**

– Mr President, firstly, well done to Mohammed, the rapporteur for the carbon border adjustment mechanism. I thought the Commission's proposals lacked ambition in terms of its timelines for phasing out free allowances and in terms of its scope. It also lacked fairness in its treatment of the least—developed countries. Mohammed has done well to address some of these problems. Least—developed countries are the least responsible for the climate crisis, yet the worst impacted by it. The two most widely discussed options to protect producers in these countries from the effects of the regulation are exemptions from the scope of the regulation and recycling of revenues generated by the regulation in order to help producers to decarbonise. I did find it incredible that the Commission's legislative proposal included neither option. Exemptions on their own would risk leaving them stranded with carbon—intensive production methods. The use of CBAM revenues therefore is crucial to allow LDCs to decarbonise their supply chains. The Commission's proposal suggests that the idea of just transition should stop at the EU borders, which is disappointing.

**Gemini 2.5 Reasoning**

The MEP's speech, while critical of the Commission's initial proposal for CBAM regarding ambition and fairness towards Least Developed Countries (LDCs), explicitly praises the rapporteur for addressing some of these problems. The MEP advocates for specific improvements (like recycling revenues for LDC decarbonisation) within the CBAM framework, indicating support for the mechanism itself provided key concerns are addressed. This suggests the MEP is likely to vote "For" the proposed legislation, assuming the rapporteur's efforts or subsequent amendments have incorporated improvements aligned with their views, particularly on LDC support, as part of the overall "Fit for 55" package.

**Gemini 2.5 Gender Prediction Error Example**

**Example: *European Central Bank – annual report 2020 (debate)***

**Ground Truth:** FEMALE          **Gemini 2.5 Prediction:** MALE (Confidence: 5)

**Gemini 2.5 Reasoning**

The speaker uses the pronoun 'Nous' followed by the masculine plural adjective 'heureux' ('Nous sommes heureux de constater...') and the masculine plural pronoun 'tous' ('Nous constatons tous...'). In French, when 'Nous' refers to a group that includes the speaker, the grammatical agreement of adjectives or certain pronouns typically reflects the speaker's gender, especially in formal speech, unless the group is entirely female (in which case the feminine form would be used). The use of 'heureux' (masculine plural) and 'tous' (masculine plural) strongly indicates that the speaker identifies as male, as a male speaker would use the masculine form regardless of the group composition, while a female speaker might use the feminine plural 'heureuses' and 'toutes' even if the group is mixed, or would use the masculine only if referring strictly to the group and not implicitly including herself in the gender agreement. The consistent use of masculine plural forms is a very strong indicator of a male speaker in this context.

**Mistral Vote Prediction Error Example**

**Example:** *Common agricultural policy - support for strategic plans to be drawn up by Member States and financed by the EAGF and by the EAFRD (debate)*

**Ground Truth:** FOR                    **Mistral Prediction:** AGAINST (Confidence: 4)

**MEP Speech**

– Monsieur le Président, agriculteurs et citoyens ont des attentes fortes vis-à-vis de la prochaine politique agricole commune. Je regrette un budget en baisse de 30 milliards d'euros environ par rapport au programme de 2014-2020, alors que le secteur agricole doit faire face à de nouveaux défis. Rappelons tout de même que la France va verser dès l'année prochaine 26,8 milliards d'euros à l'Union européenne – c'est un record nous concernant – et que nous avons besoin de la PAC. Un débat public est actuellement en cours sur les enjeux de la prochaine PAC, dans une relative indifférence médiatique. Pourtant, les propositions pourraient être reprises dans le plan stratégique national qui sera la déclinaison française de la PAC, replaçant ainsi, comme il se devrait, les États à l'initiative des décisions. Notre modèle agricole doit désormais affronter le défi écologique de la transition verte, mais sa mise en œuvre est à affiner entre la nouvelle structure verte de la PAC et les mesures environnementales inscrites dans la stratégie F2F. Les accords commerciaux de l'Union européenne, et notamment l'accord unique de libre-échange UE-Mercosur, suscitent des inquiétudes légitimes sur la protection de nos normes de qualité dans un secteur déjà fortement impacté par la pandémie et la sécheresse. Gageons que la réforme prévue au premier trimestre 2021 permette de développer les circuits courts et un revenu digne pour nos agriculteurs.

**Mistral Gender Prediction Error Example**

**Example: *Guidelines for the 2024 budget - Section III (debate)***

**Ground Truth:** MALE                    **Mistral Prediction:** FEMALE (Confidence: 5)

**MEP Speech**

I have been involved in the works of the budgetary procedure for a few years, and now I must say that from the position of the FEMM Committee, it resembles 'Groundhog Day'. We repeatedly called for the Citizens, Equality, Rights and Values Programme to include separate gender-related budget lines, with no success. We called for the Daphne programme resources to increase the special budget for gender-related issues, but with no success. We continued to call for the connection between attacks on the rule of law and worrisome backlash against gender equality and women's rights in the EU, and attacks on LGBTI+ people will be connected to the budget issues, without success. And what's most important, we repeatedly call for gender equality perspective to be fully integrated into all policy areas, at all stages of the budgetary process. Yet, the European Court of Auditors found again that gender mainstreaming is the least successfully implemented overarching policy priority, again and again and again, without success.

**Mistral Reasoning**

The speech is entirely focused on gender equality, women's rights, LGBTI+ rights, and budgetary procedures related to these issues. The reference to speaking "from the position of the FEMM Committee" (Committee on Women's Rights and Gender Equality) and the repeated emphasis on topics traditionally championed by women led the model to infer a female identity. While male MEPs can certainly champion these causes, the intensity of advocacy and association with the FEMM agenda resulted in a high-confidence misclassification.

**Claude 3.5 Vote Prediction Error Example**

**Example: *2019–2020 Reports on Turkey (debate)***

**Ground Truth:** AGAINST                    **Claude 3.5 Prediction:** FOR (Confidence: 5)

### MEP Speech

– Voorzitter, ik kan mijn ogen niet geloven. Bladzijden lang wordt beschreven hoe Turkije nog steeds alle afspraken aan zijn militaire laarzen lapt. Er zijn duizenden politiek gevangenen, onder wie zwangere vrouwen, er wordt gefolterd, mensen verdwijnen, 152 000 rechters, leraren, journalisten werden ontslagen. De Turken bezetten al een halve eeuw ongestraft een deel van onze lidstaat Cyprus. Zij chanteren ons met een massale vluchtelingenstroom. En toch, toch omschrijft men vandaag Turkije, het islamitische Turkije, hier als een strategische partner op geopolitiek en economisch vlak, die bovendien lid moet blijven van de NAVO. Na het recente bombarderen van Koerden en Armeniërs zou men denken dat de eurobonzen beseffen dat Turkije een oorlogsstoker is en niet een partner om een stabiel en veilig klimaat mee te creëren in de oostelijke Middellandse Zee. Hoe naïef kan men zijn?

### Claude 3.5 Reasoning

The MEP's speech is overwhelmingly critical of Turkey, listing numerous severe human rights abuses, undemocratic actions, and foreign policy violations. The use of emotive and condemnatory language suggests strong disapproval. Claude 3.5 interpreted this negativity as support for a critical report, hence predicting a 'FOR' vote. However, the actual vote was 'AGAINST', potentially reflecting opposition to the EU's framing or political motives behind the report, not the content itself.

**Claude 3.5 Gender Prediction Error Example**

### Example: *Conclusion of the EU-Viet Nam Free Trade Agreement (debate)*

**Ground Truth:** FEMALE          **Claude 3.5 Prediction:** MALE (Confidence: 4)

### MEP Speech

Sans discernement dans l'octroi des préférences généralisées, dont le Viêt Nam bénéficie, elle aura été tout aussi inepte dans les négociations commerciales. Les échanges avec ce pays devaient pourtant la placer en position de force. Il n'en est rien. Un exemple cocasse le démontre: alors que l'Union européenne reconnaît 39 IGP vietnamiennes, le Viêt Nam n'en reconnaîtrait que 36 françaises, dont le champagne – avec une période transitoire de dix ans!
L'accord n'est évidemment pas une bonne nouvelle en matière de droits de l'homme. Outre la liberté syndicale, en 2012, 1,75 million de Vietnamiens âgés de 5 à 17 ans étaient toujours concernés par le travail des enfants. Quant à nos industries – notamment le textile, historiquement implanté dans le Grand Est – elles paieront par une nouvelle saignée cette course au moins-disant social, les salaires vietnamiens plafonnant à 59 % des salaires chinois, ce qui explique vraisemblablement l'empressement déployé pour conclure cet accord.

### Claude 3.5 Reasoning

The text employs a formal, assertive, and analytical tone, focusing on technical and economic details such as trade agreements, labor statistics, and industry impacts. This style is often associated with traditionally male communication patterns in political and economic discourse. The use of precise figures, structured argumentation, and a critical stance also aligns with stereotypical male linguistic tendencies in professional contexts. However, as gendered language patterns are not definitive, the confidence level is set at 4.

# Assessing ASR Robustness for Burmese: Impacts of Missing Speech Segments and Interruptions

**Ankit Maurya[1]    Manikandan Ravikiran[1]    Rohit Saluja[1]**
[1]Indian Institute of Technology Mandi
{s24090,erpd2301}@students.iitmandi.ac.in, rohit@iitmandi.ac.in

## Abstract

The study explores the performance, robustness, and effects of automatic speech recognition systems when speech is missing or interrupted, with a specific focus on Burmese, a low-resource language. This study addresses several key research questions: How does missing or interrupted speech affect the accuracy of ASR? What is the link between the length of missing speech and the accuracy of the transcription? How are errors propagating when speech is masked or interrupted? By fine-tuning Wav2vec-bert2.0 and MMS-Zeroshot-300M (Massively Multilingual Speech) on a regular speech dataset (OpenSLR) of Burmese, the study answers these questions by evaluating the models on OpenSLR and 2 other datasets (FLEURS and Bloom) on common ASR metrics like Word Error Rate and Character Error Rate. The results reveal significant insights into error propagation, positional error patterns, and dataset-specific robustness. The study provides a baseline and methodological insights for future ASR research in interrupted settings for low-resource languages. The study's findings can inform the development of more robust ASR systems for real-world applications in low-resource languages.

## 1 Introduction

While automatic speech recognition (ASR) has enabled applications ranging from voice assistants (Dubiel et al., 2018; Sim et al., 2019) to automated transcription services (Jeffries et al., 2024), its performance often degrades under real-world conditions involving missing or incomplete speech segments (Barker et al., 2013; Gemmeke et al., 2011). Missing speech in these scenarios commonly arises from network packet loss during VoIP or streaming (Dissen et al., 2024; Lee and Kang, 2013; Kumalija and Nakamoto, 2022), recording interruptions caused by hardware issues or user-generated noise, and transmission errors from corrupted media. This challenge is particularly acute

for low-resource languages (LRLs), which have not seen the same focus on robustness as high-resource languages (Baevski et al., 2020; Rubenstein et al., 2023; Radford et al., 2022), raising important questions about ASR reliability in these contexts.

Accordingly, this study focuses on Burmese, a Tibeto-Burman language spoken by approximately 42.9 million people and notably underrepresented in ASR research (Wikipedia contributors, 2025b; Li and Jian, 2024). The linguistic complexity of Burmese makes it a compelling case for robustness analysis. As a tonal language with an agglutinative morphology (Wikipedia contributors, 2025a), meaning is conveyed through subtle changes in pitch and duration, while grammatical information is often encoded in extended word forms. Consequently, the limited data available for training robust models poses a critical challenge, as even brief gaps in audio can lead to significant loss of semantic and grammatical information. Enhancing ASR robustness for Burmese can thus bridge communication gaps and enable technological inclusion for millions of speakers. To investigate these challenges, this paper aims to:

- Quantify the performance impact of missing speech on a fine-tuned Burmese ASR model;

- Perform a detailed error analysis to understand the nature and distribution of errors caused by missing data; and

- Evaluate model robustness across multiple relevant datasets.

The remainder of this paper is organized as follows. Section 2 reviews related work. Section 3 describes our experimental setup, including the datasets and methodology. Finally, Section 5 concludes the paper and outlines future work.

## 2 Related Works

The challenge of handling missing audio in automatic speech recognition (ASR) has long been a tough problem. In the past, researchers mostly focused on techniques like robust feature extraction and model adaptation to reduce the effects of noise and distortions (Ming and Crookes, 2014). More recently, the focus has shifted to what's called speech inpainting or speech reconstruction, which involves filling in or estimating missing segments of audio. These methods range from simple interpolation approaches (Kauppinen et al., 2001) to more advanced models like Gaussian Mixture Models (GMMs) (Cooke et al., 2001) and deep learning techniques using autoencoders or Generative Adversarial Networks (GANs) (Wali et al., 2021).

Another related idea comes from self-supervised learning, where models are trained to predict masked parts of the input data. This was first popularized in natural language processing by models like BERT (Devlin et al., 2019) and later adapted for speech in models like Wav2vec 2.0 (Baevski et al., 2020) and HuBERT (Hsu et al., 2021). These models learn robust representations that help in handling missing data. Though primarily used for pre-training, masking is also used as a data augmentation strategy to make models more robust (Rebuffi et al., 2021). One recent example is SpeechPainter (Zalan Borsos and Matthew Sharifi and Marco Tagliasacchi, 2022), which uses a generative diffusion model to do high-quality speech inpainting.

However, most of these efforts have focused on high-resource languages. In our study, we explore how masking and gapping two strategies for simulating missing speech, differently affect ASR error patterns and robustness, specifically for Burmese.

## 3 Experimental Setup

This section outlines the experimental framework designed to assess the performance and robustness of ASR systems for Burmese with missing speech.

### 3.1 Dataset

To evaluate the impact of missing speech on Burmese ASR, we utilize publicly available Burmese speech datasets, representing different domains, recording conditions, and potentially varying levels of annotation quality. The primary datasets considered are: (a) **OpenSLR Burmese**

**(SLR80)** (Oo et al., 2020): Contains approximately 4 hours of read speech, used here for fine-tuning and evaluation. Characteristics often include relatively clean recordings. (b) **Few-shot Learning Evaluation of Universal Representations of Speech (FLEURS)** (Conneau et al., 2023; Goyal et al., 2022):An n-way parallel dataset with 12 hours of speech per language, representing a more diverse source. (c) **Bloom Dataset** (Leong et al., 2022): Contains 1 hour of burmese read speech from book paragraphs. Table 1 presents the dataset splits we use in our study.

| Dataset | Split Used | No. of Samples |
|---------|-----------|----------------|
| OpenSLR | Test + Val | 196+206 |
| FLEURS | Test | 880 |
| BLOOM | Test + Val | 50+50 |

Table 1: Dataset and Split Configuration for the Study.

To analyze ASR model robustness to incomplete speech, we simulate missing data using two controlled techniques: masking and gapping. These are applied over specified durations $D_m \in \{0.25\text{s}, 0.5\text{s}, 0.75\text{s}, 1\text{s}\}$, chosen to reflect real-world interruption lengths, such as network packet loss in VoIP systems (e.g., 10–40 packets of 25ms) (Zhang et al., 2024). Let an original discrete-time audio signal be $S_{orig}[n]$, with $N_{orig}$ samples. The gap/mask duration $D_m$ corresponds to $N_D$ samples. If $N_{orig}$ is not a multiple of $N_D$, $S_{orig}[n]$ is zero-padded at the end to form a signal $S[n]$ of length $N$, where $N$ is the smallest multiple of $N_D$ such that $N \geq N_{orig}$. Subsequent operations refer to this signal $S[n]$ of length $N$.

**Masking** Masking simulates data loss by setting a segment of $N_D$ samples to zero amplitude, with the masked signal $S'_{mask,k}[n]$ retaining the length $N$. For each signal $S[n]$, $N/N_D$ distinct masked versions are generated. The $k^{th}$ masked signal is formed by zeroing out samples in the segment $n \in [kN_D, (k+1)N_D - 1]$, where $k \in [0, \ldots, (N/N_D) - 1]$. This is described by:

$$S'_{mask,k}[n] = \begin{cases} S[n] & \text{for } 0 \leq n < kN_D \\ 0 & \text{for } kN_D \leq n < (k+1)N_D \\ S[n] & \text{for } (k+1)N_D \leq n < N \end{cases}$$

This technique is analogous to how lost data packets replaced by silence are handled and is fundamental in applying time-frequency masks in speech enhancement (Kim, 2021). Alternatively, it is an element-wise multiplication $S'_{mask,k} = S \odot M_k$, where $M_k[n]$ is a binary mask.

441

**Gapping**  Gapping simulates interruptions by inserting a silent segment of $N_D$ zero-valued samples into $S[n]$, thereby extending its total duration to $N + N_D$ samples. For each signal $S[n]$, $(N/N_D) + 1$ distinct gapped versions are generated. In the $j^{th}$ gapped signal, $S'_{gap,j}[m]$, silence is inserted at a position corresponding to $jN_D$ in the timeline of $S[n]$, where $j \in [0, \ldots, N/N_D]$. The resulting signal $S'_{gap,j}[m]$ for $m \in [0, N + N_D - 1]$ is constructed as:

$$S'_{gap,j}[m] = \begin{cases} S[m], & \text{for } 0 \le m < jN_D, \\ 0, & \text{for } jN_D \le m < jN_D + N_D, \\ S[m - N_D], & \text{for } jN_D + N_D \le m < N + N_D \end{cases}$$

This method alters the signal's temporal structure, testing the model's ability to handle unexpected pauses. The key distinction is that masking overwrites existing audio data while preserving effective signal duration ($N$), whereas gapping inserts new silent data, increasing the overall signal length and specifically challenging resilience to pauses and temporal shifts.

### 3.2  Models and Parameters

In this work, we selected Wav2Vec-BERT-2.0 (Chung et al., 2021) and MMS-Zeroshot-300M [1] (Pratap et al., 2024) due to their extensive multilingual pre-training and strong performance on LRLs. Wav2Vec-BERT-2.0's prior exposure to Burmese and MMS-Zeroshot-300M's training on over 1,000 languages make them suitable for this study. We fine-tune publicly available checkpoints for our experiments. This study fine-tunes both models on the OpenSLR dataset for up to 16 epochs using the AdamW optimizer (learning rate 5e-5, batch size 8), utilizing pretrained checkpoints available on Hugging Face  (Face, 2025) and utilizing the Transformers library for the fine-tuning process. For our experiments, we select the best-performing fine-tuned checkpoints of both Wav2Vec-BERT-2.0 and MMS-Zeroshot-300M on the OpenSLR validation set with the best validation WER and use it for all subsequent experiments. All subsequent evaluations use the data splits detailed in Table 1. All the experiments were done on a single NVIDIA RTX A5000 GPU.

### 3.3  Evaluation Metrices

We assess ASR performance using the following metrics to quantify transcription errors and to analyze the impact of missing speech:

- **WER (Word Error Rate)**: Measures word-level transcription errors.

- **CER (Character Error Rate)**: Measures character-level transcription errors.

- **Error Percentage (Error %)**: The proportion of samples where the transcription differs from the clean audio baseline after simulating missing speech.

- **Edit Distance Distribution**: Examines how severe the errors are by analyzing the distribution of edit distances for each affected sample compared to the baseline.

- **Positional Error Analysis**: Looks at where errors tend to occur, focusing on the first and last segments around the missing speech region.

| Dur. | Masking | | Gapping | |
|------|---------|---------|---------|---------|
| | Samples | Err. (%) | Samples | Err. (%) |
| 0.25s | 9144 | 69.14 | 9546 | 77.92 |
| 0.50s | 4458 | 84.28 | 4860 | 77.72 |
| 0.75s | 2924 | 88.68 | 3326 | 78.35 |
| 1s | 2134 | 91.38 | 2536 | 74.80 |

Table 2: Total processed samples and error percentage (%) for Masking vs. Gapping on the OpenSLR dataset by duration. (Model Used: Wav2Vec-BERT-2.0)

## 4  Results and Analysis

Across all three datasets, OpenSLR, Fleurs, and Bloom, our experiments highlight that the impact of missing speech (via masking and gapping) is closely tied to baseline dataset characteristics and duration of gap/mask. Trends across datasets are captured in Appendix Figures 7, 8 (overall WER/CER), Figure 1 (OpenSLR edit distance), Figures 4, 5 (Fleurs), Figures 3, 6 (Bloom), and Tables 3 and 4.

**OpenSLR (Appendix Figures 1, 2, 7):** This relatively clean dataset exhibits clear, progressive increases in both WER and CER as mask duration grows. Masking is notably more detrimental than gapping, with Wav2Vec-BERT-2.0's CER rising from 6.79% to 21.02% at 1s (compared to 8.89% for gapping). Edit distance distributions (Figure 1) confirm more severe character-level errors for longer mask durations. Positional analysis

(Figure 2) reveals that masking the final segment of an utterance dramatically increases edit distances compared to initial-segment masking, underscoring the importance of utterance-end information for ASR. This can be attributed to the model's need for forward and backward context to correctly decode speech. When the end of an utterance is masked, the model loses critical cues for disambiguation, which is especially important for a tonal and agglutinative language like Burmese where meaning and grammatical information can be conveyed by subtle changes at the end of words. The absence of this final information results in a significantly higher character-level error rate. Error percentages (Table **??**) mirror these patterns, especially for masking (69.14% to 91.38% with duration).

**Fleurs (Appendix Figures 4, 5 ):** Fleurs is substantially more challenging, with baseline WERs already above 179% and CERs over 21% even without missing speech. Masking consistently elevates CER (Wav2Vec-BERT-2.0: 21.26% to 25.35%), while WER shows only minor, sometimes negative shifts likely due to error saturation. Edit distance distributions (Figure 4) and average edit distances (Figure 5) confirm more fine-grained character errors from masking, especially at utterance ends. Gapping causes a severe but consistent level of degradation across all gap durations, with WER and CER remaining largely stable ( Wav2Vec-BERT-2.0 WER ≈184%, CER ≈22%). Error percentages (Table 4) remain extremely high (>90%) across all gap/mask types and durations.

**Bloom (Appendix Figures 3, 6):** The Bloom dataset poses the greatest challenge, with baseline CERs over 66% and WERs exceeding 150%. Neither masking nor gapping meaningfully alters the already saturated error rates (CER ≈67-68%, WER ≈150-160%). Edit distance distributions and error percentages (near 100% for all conditions) confirm that the models nearly uniformly fail on this dataset, irrespective of the missing speech scenario. Overall, our key takeaways include:

- **Masking** generally introduces more severe character-level errors (as seen in CER and edit distance shifts), particularly on cleaner datasets like OpenSLR and moderately so on Fleurs. This highlights that real-world ASR systems deployed in environments with short-duration or partial occlusions (e.g., coughs, short microphone dropouts) are likely to see disproportionately larger transcription errors, especially on simpler, cleaner audio inputs.

- **Gapping** causes more stable but consistently high errors, with less sensitivity to gap/mask duration. This suggests that silent gaps or short audio losses (e.g., packet loss in VoIP or poor connectivity) might degrade performance consistently across a range of scenarios, rather than in a duration-dependent manner.

- **Dataset Difficulty Dominates:** On challenging datasets (Fleurs, Bloom), extremely high baseline errors overshadow the incremental effects of missing speech, leading to error saturation. This indicates that for real-world ASR robustness, improving baseline model performance (e.g., adapting to domain-specific vocabularies, reducing dataset-domain mismatch) is critical, as gap/mask effects become secondary when baseline transcription itself is unreliable.

These insights reinforce that efforts to improve ASR robustness in real-world scenarios must prioritize both baseline domain adaptation and gap or mask specific resilience, especially for character-level fidelity and utterance-end information crucial for downstream tasks.

## 5 Conclusion

In conclusion, this study reveals significant challenges in applying any ASR system to LRLs when encountering simulated missing speech segments. Our findings indicate a clear degradation in performance, measured by WER and CER, as the duration or ratio of missing data increases. Notably, the structure of the missing data matters; our results suggest that masking the audio signal generally causes more severe degradation than inserting gaps. Furthermore, the errors induced by these missing segments are not localized; they often propagate beyond the immediate vicinity of the gap or mask, primarily manifesting as deletions and substitutions, which aligns with observations from edit distance analysis. Finally, the system's resilience is not uniform, varying significantly across the different datasets tested (OpenSLR, Fleurs, Bloom), underscoring the influence of acoustic conditions, speaking styles, and domain specificity on robustness. These results highlight the need for targeted strategies to improve the robustness of Burmese ASR

systems against various forms of missing speech data encountered in real-world scenarios.

## Limitations

This study's limitations include the simulation of missing speech (zero-masking, gapping), which may not fully mirror real-world data loss complexities. Findings are also specific to the chosen pre-trained models (Wav2Vec-BERT-2.0, MMS-Zeroshot-300M) potentially differing for other ASR architectures or training methods. Furthermore, the employed Burmese datasets, while valuable, may not encompass the language's complete dialectal or acoustic diversity, which could affect the broader generalizability of the observed robustness levels.

## Future Works

Future work will focus on three key areas. First, we will explore advanced speech inpainting and reconstruction techniques to better mitigate the effects of missing speech. Second, we plan to investigate more dynamic and realistic interruption patterns, such as randomly distributed or non-uniform segment lengths, to provide a more accurate assessment of ASR robustness. Finally, we will consider a broader range of Burmese datasets to improve the generalizability of our findings. We also plan to evaluate models with shorter durations, such as 0.10s, to capture more nuanced effects on performance.

## Acknowledgments

## Ethics Statement

This research exclusively utilizes publicly available speech datasets (OpenSLR, Fleurs, and Bloom) intended for academic use, and no new data was collected from human subjects. The primary goal of this work is to foster positive technological inclusion by analyzing and improving ASR robustness for the low-resource Burmese-speaking community. The findings are specific to the models and datasets employed, which may not encompass the full linguistic and dialectal diversity of the Burmese language, potentially leading to performance disparities across different speaker populations. The computational cost required to train and run these large models also represents an environmental consideration.

## References

Alexei Baevski, Yuhao Zhou, Abdelrahman Mohamed, and Michael Auli. 2020. wav2vec 2.0: A framework for self-supervised learning of speech representations. In *Advances in Neural Information Processing Systems*, volume 33, pages 12449–12460. Curran Associates, Inc.

Jon Barker, Emmanuel Vincent, Ning Ma, Heidi Christensen, and Phil Green. 2013. The pascal chime speech separation and recognition challenge. *Computer Speech Language*, pages 621–633.

Yu-An Chung, Yu Zhang, Wei Han, Chung-Cheng Chiu, James Qin, Ruoming Pang, and Yonghui Wu. 2021. W2v-bert: Combining contrastive learning and masked language modeling for self-supervised speech pre-training. In *2021 IEEE Automatic Speech Recognition and Understanding Workshop (ASRU)*, pages 244–250. IEEE.

Alexis Conneau, Min Ma, Simran Khanuja, Yu Zhang, Vera Axelrod, Siddharth Dalmia, Jason Riesa, Clara Rivera, and Ankur Bapna. 2023. Fleurs: Few-shot learning evaluation of universal representations of speech. In *2022 IEEE Spoken Language Technology Workshop (SLT)*, pages 798–805.

Martin Cooke, Phil Green, Ljubomir Josifovski, and Ascension Vizinho. 2001. Robust automatic speech recognition with missing and unreliable acoustic data. *Speech Communication*, 34(3):267–285.

Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. 2019. Bert: Pre-training of deep bidirectional transformers for language understanding. In *North American Chapter of the Association for Computational Linguistics*.

Yehoshua Dissen, Shiry Yonash, Israel Cohen, and Joseph Keshet. 2024. Enhanced asr robustness to packet loss with a front-end adaptation network. *ArXiv*, abs/2406.18928.

Mateusz Dubiel, Martin Halvey, and Leif Azzopardi. 2018. A survey investigating usage of virtual personal assistants. *ArXiv*, abs/1807.04606.

Hugging Face. 2025. Wav2vec2-bert model documentation. https://huggingface.co/docs/transformers/en/model_doc/wav2vec2-bert. Accessed: 2025-04-25.

Jort F. Gemmeke, Maarten Van Segbroeck, Yujun Wang, Bert Cranen, and Hugo Van Hamme. 2011. Automatic speech recognition using missing data techniques: Handling of real-world data. In *Robust Speech Recognition of Uncertain or Missing Data*.

Naman Goyal, Cynthia Gao, Vishrav Chaudhary, Peng-Jen Chen, Guillaume Wenzek, Da Ju, Sanjana Krishnan, Marc'Aurelio Ranzato, Francisco Guzmán, and Angela Fan. 2022. The Flores-101 evaluation benchmark for low-resource and multilingual machine translation. *Transactions of the Association for Computational Linguistics*, 10:522–538.

Wei-Ning Hsu, Benjamin Bolte, Yao-Hung Hubert Tsai, Kushal Lakhotia, Ruslan Salakhutdinov, and Abdel-rahman Mohamed. 2021. Hubert: Self-supervised speech representation learning by masked prediction of hidden units.

Nat Jeffries, Evan King, Manjunath Kudlur, Guy Nicholson, James Wang, and Pete Warden. 2024. Moonshine: Speech recognition for live transcription and voice commands.

I. Kauppinen, Jyrki K. Kauppinen, and Pekka Saarinen. 2001. A method for long extrapolation of audio signals. *Journal of The Audio Engineering Society*, 49:1167–1180.

Gibak Kim. 2021. Review of time–frequency masking approach for improving speech intelligibility in noise. *IETE Technical Review*, 39:1–12.

Elhard Kumalija and Yukikazu Nakamoto. 2022. Performance evaluation of automatic speech recognition systems on integrated noise-network distorted speech. *Frontiers in Signal Processing*, 2:999457.

Min-Ki Lee and Hong-Goo Kang. 2013. Speech quality estimation of voice over internet protocol codec using a packet loss impairment model. *The Journal of the Acoustical Society of America*, 134:EL438–44.

Colin Leong, Joshua Nemecek, Jacob Mansdorfer, Anna Filighera, Abraham Owodunni, and Daniel Whitenack. 2022. Bloom library: Multimodal datasets in 300+ languages for a variety of downstream tasks. In *Proceedings of the 2022 Conference on Empirical Methods in Natural Language Processing*, pages 8608–8621, Abu Dhabi, United Arab Emirates. Association for Computational Linguistics.

Yaxuan Li and Yang Jian. 2024. Low-resource burmese speech synthesis based on visual text embedding and diffusion model. In *Proceedings of the 2024 4th International Conference on Artificial Intelligence, Automation and Algorithms*, AI2A '24, page 34–40, New York, NY, USA. Association for Computing Machinery.

Ji Ming and Danny Crookes. 2014. Speech enhancement from additive noise and channel distortion - a corpus-based approach. In *Interspeech*.

Yin May Oo, Theeraphol Wattanavekin, Chenfang Li, Pasindu De Silva, Supheakmungkol Sarin, Knot Pipatsrisawat, Martin Jansche, Oddur Kjartansson, and Alexander Gutkin. 2020. Burmese Speech Corpus, Finite-State Text Normalization and Pronunciation Grammars with an Application to Text-to-Speech. In *Proceedings of The 12th Language Resources and Evaluation Conference (LREC)*, pages 6328–6339, Marseille, France. European Language Resources Association (ELRA).

Vineel Pratap, Andros Tjandra, Bowen Shi, Paden Tomasello, Arun Babu, Sayani Kundu, Ali Elkahky, Zhaoheng Ni, Apoorv Vyas, Maryam Fazel-Zarandi, et al. 2024. Scaling speech technology to 1,000+ languages. *Journal of Machine Learning Research*, 25(97):1–52.

Alec Radford, Jong Wook Kim, Tao Xu, Greg Brockman, Christine McLeavey, and Ilya Sutskever. 2022. Robust speech recognition via large-scale weak supervision. In *International Conference on Machine Learning*.

Sylvestre-Alvise Rebuffi, Sven Gowal, Dan Andrei Calian, Florian Stimberg, Olivia Wiles, and Timothy Mann. 2021. Data augmentation can improve robustness. In *Neural Information Processing Systems*.

Paul K. Rubenstein, Chulayuth Asawaroengchai, Duc Dung Nguyen, Ankur Bapna, Zalán Borsos, Félix de Chaumont Quitry, Peter Chen, Dalia El Badawy, Wei Han, Eugene Kharitonov, Hannah Muckenhirn, Dirk Padfield, James Qin, Danny Rozenberg, Tara Sainath, Johan Schalkwyk, Matt Sharifi, Michelle Tadmor Ramanovich, Marco Tagliasacchi, Alexandru Tudor, Mihajlo Velimirović, Damien Vincent, Jiahui Yu, Yongqiang Wang, Vicky Zayats, Neil Zeghidour, Yu Zhang, Zhishuai Zhang, Lukas Zilka, and Christian Frank. 2023. Audiopalm: A large language model that can speak and listen.

Khe Chai Sim, Petr Zadrazil, and Françoise Beaufays. 2019. An investigation into on-device personalization of end-to-end automatic speech recognition models. In *Interspeech 2019*. ISCA.

Aamir Wali, Zareen Alamgir, Saira Karim, Ather Fawaz, Mubariz Ali, Muhammad Adan, and Malik Mujtaba. 2021. Generative adversarial networks for speech processing: A review. *Computer Speech & Language*, 72:101308.

Wikipedia contributors. 2025a. Burmese grammar — Wikipedia, the free encyclopedia. [Online; accessed 18-May-2025].

Wikipedia contributors. 2025b. Burmese language — Wikipedia, the free encyclopedia. [Online; accessed 2-April-2025].

Zalan Borsos and Matthew Sharifi and Marco Tagliasacchi. 2022. SpeechPainter: Text-conditioned Speech Inpainting. In *Interspeech 2022*, pages 431–435.

Zihan Zhang, Jiayao Sun, Xianjun Xia, Chuanzeng Huang, Yijian Xiao, and Lei Xie. 2024. Bs-plcnet: Band-split packet loss concealment network with multi-task learning framework and multi-discriminators. In *2024 IEEE International Conference on Acoustics, Speech, and Signal Processing Workshops (ICASSPW)*, pages 23–24. IEEE.

# A Appendix



Figure 1: Edit distance distribution across samples for OpenSLR dataset. (Model: Wav2Vec-BERT-2.0)



Figure 2: Average Edit distance for First and Last segment of mask and gap audio segments. (Model= Wav2Vec-BERT-2.0, Dataset = OpenSLR) (Baseline = Prediction without mask/gap)



Figure 3: Average Edit distances for First segment and Last Segment. (Model = Wav2Vec-BERT-2.0, dataset = Bloom)



Figure 4: Edit Distance distribution of Fleurs dataset. (Model: Wav2Vec-BERT-2.0)



Figure 5: Average Edit distances for First segment and Last Segment. (Model = Wav2Vec-BERT-2.0, dataset = Fleurs)



Figure 6: Edit Distance distribution of Bloom dataset. (model = Wav2Vec-BERT-2.0)

# CER Trends Across Models and Datasets



Figure 7: CER Trend for all datasets for different gap and mask duration for both models.

Table 3: Detailed WER (%) and CER (%) Results by Duration for Masking vs. Gapping Across Datasets and Models. Models were fine-tuned (FT) on OpenSLR.

| Dataset | Model | Duration of mask/gap | Masking | | Gapping | |
|---------|-------|---------------------|---------|---------|---------|---------|
| | | | WER (%) | CER (%) | WER (%) | CER (%) |
| OpenSLR | Wav2Vec-BERT-2.0 (FT) | 0s | 37.35 | 6.79 | 37.35 | 6.79 |
| | | 0.25s | 40.85 | 8.60 | 41.39 | 9.06 |
| | | 0.50s | 44.75 | 12.70 | 41.29 | 9.28 |
| | | 0.75s | 47.69 | 16.82 | 41.53 | 9.09 |
| | | 1s | 50.46 | 21.02 | 41.12 | 8.89 |
| | MMS-Zeroshot-300M (FT) | 0s | 56.37 | 12.03 | 56.37 | 12.03 |
| | | 0.25s | 59.94 | 14.50 | 59.78 | 14.40 |
| | | 0.50s | 61.74 | 18.33 | 61.29 | 17.86 |
| | | 0.75s | 63.81 | 22.26 | 63.00 | 21.17 |
| | | 1s | 65.58 | 26.10 | 64.38 | 24.31 |
| Fleurs | Wav2Vec-BERT-2.0 (FT) | 0s | 181.00 | 21.26 | 181.00 | 21.26 |
| | | 0.25s | 182.60 | 22.13 | 184.00 | 22.34 |
| | | 0.50s | 180.10 | 23.18 | 184.00 | 22.39 |
| | | 0.75s | 177.00 | 24.24 | 184.00 | 22.34 |
| | | 1s | 174.40 | 25.35 | 184.00 | 22.29 |
| | MMS-Zeroshot-300M (FT) | 0s | 179.24 | 28.49 | 179.24 | 28.49 |
| | | 0.25s | 178.50 | 29.51 | 180.29 | 29.38 |
| | | 0.50s | 175.42 | 30.48 | 179.90 | 29.36 |
| | | 0.75s | 172.34 | 31.50 | 179.48 | 29.41 |
| | | 1s | 169.05 | 32.59 | 178.95 | 29.50 |
| Bloom | Wav2Vec-BERT-2.0 (FT) | 0s | 182.45 | 66.90 | 182.45 | 66.90 |
| | | 0.25s | 183.66 | 66.65 | 186.10 | 66.98 |
| | | 0.50s | 180.94 | 66.75 | 185.32 | 66.98 |
| | | 0.75s | 178.19 | 67.09 | 185.22 | 67.13 |
| | | 1s | 175.10 | 67.47 | 186.24 | 67.08 |
| | MMS-Zeroshot-300M (FT) | 0s | 159.27 | 67.99 | 159.27 | 67.99 |
| | | 0.25s | 157.45 | 67.76 | 159.27 | 68.01 |
| | | 0.50s | 154.67 | 67.93 | 159.21 | 67.85 |
| | | 0.75s | 151.84 | 68.28 | 159.01 | 67.60 |
| | | 1s | 148.55 | 68.66 | 158.13 | 67.87 |

## WER Trends Across Models and Datasets



Figure 8: WER Trend for all datasets for different gap and mask durations for both models.

Table 4: Total Processed Samples and Percentage of Samples with Errors/Prediction Changes Across Datasets, Models, and Conditions. Models were fine-tuned (FT) on OpenSLR.

| Dataset | Model | Duration of mask/gap | Masking | | Gapping | |
|---|---|---|---|---|---|---|
| | | | Total Samples | Error (%) | Total Samples | Error (%) |
| OpenSLR | Wav2Vec-BERT-2.0 (FT) | 0.25s | 9144 | 69.14 | 9546 | 77.92 |
| | | 0.50s | 4458 | 84.28 | 4860 | 77.72 |
| | | 0.75s | 2924 | 88.68 | 3326 | 78.35 |
| | | 1s | 2134 | 91.38 | 2536 | 74.80 |
| | MMS-Zeroshot-300M (FT) | 0.25s | 9144 | 79.01 | 9546 | 76.79 |
| | | 0.50s | 4458 | 87.01 | 4860 | 82.88 |
| | | 0.75s | 2924 | 90.94 | 3326 | 84.97 |
| | | 1s | 2134 | 92.83 | 2536 | 86.24 |
| Fleurs | Wav2Vec-BERT-2.0 (FT) | 0.25s | 54 515 | 88.93 | 55 395 | 94.00 |
| | | 0.50s | 27 037 | 94.02 | 27 917 | 94.63 |
| | | 0.75s | 17 882 | 95.76 | 18 762 | 95.81 |
| | | 1s | 13 293 | 96.83 | 14 173 | 93.69 |
| | MMS-Zeroshot-300M (FT) | 0.25s | 54 515 | 82.24 | 55 395 | 91.27 |
| | | 0.50s | 27 037 | 88.77 | 27 917 | 87.45 |
| | | 0.75s | 17 882 | 91.98 | 18 762 | 94.89 |
| | | 1s | 13 293 | 94.21 | 14 173 | 92.85 |
| Bloom | Wav2Vec-BERT-2.0 (FT) | 0.25s | 3108 | 99.81 | 3208 | 99.75 |
| | | 0.50s | 1524 | 99.93 | 1625 | 99.75 |
| | | 0.75s | 1008 | 99.90 | 1108 | 99.55 |
| | | 1s | 737 | 100.00 | 838 | 99.40 |
| | MMS-Zeroshot-300M (FT) | 0.25s | 3108 | 97.52 | 3208 | 99.13 |
| | | 0.50s | 1524 | 99.61 | 1625 | 98.58 |
| | | 0.75s | 1008 | 99.90 | 1108 | 99.46 |
| | | 1s | 737 | 100.00 | 838 | 99.52 |

# Following Route Instructions using Large Vision-Language Models:
# A Comparison between Low-level and Panoramic Action Spaces

**Vebjørn Haug Kåsene**
University of Oslo
vebjorhk@gmail.com

**Pierre Lison**
Norwegian Computing Center
plison@nr.no

## Abstract

*Vision-and-Language Navigation* (VLN) refers to the task of enabling autonomous robots to navigate unfamiliar environments by following natural language instructions. While recent Large Vision-Language Models (LVLMs) have shown promise in this task, most current VLN systems rely on models specifically designed and optimized for navigation, leaving the potential of off-the-shelf LVLMs underexplored. Furthermore, while older VLN approaches used low-level action spaces with egocentric views and atomic actions (such as "turn left" or "move forward"), newer models tend to favor panoramic action spaces with discrete navigable viewpoints. This paper investigates (1) whether off-the-shelf LVLMs (fine-tuned without architectural modifications or simulator-based training) can effectively support VLN tasks and (2) whether such models can support both low-level and panoramic action paradigms. To this end, we fine-tune the open-source model Qwen2.5-VL-3B-Instruct on the *Room-to-Room (R2R)* dataset and evaluate its empirical performance across both low-level and panoramic action spaces. The best resulting model achieves a 41% success rate on the R2R test set, demonstrating that while off-the-shelf LVLMs can learn to perform Vision-and-Language Navigation, they still lag behind models specifically designed for this task.

## 1 Introduction

Mobile robots deployed in real-world environments are often tasked with reaching specific locations described in natural language. For example, a robot might be instructed to "deliver a package to the office at the end of the hallway," without prior knowledge of the environment. In such cases, a human can provide guidance through route instructions such as "Walk down the hallway and take the last door to your left." To perform its task, the robot must first interpret the linguistic input provided by the human user, ground this input in its visual perception of the environment, and execute the corresponding sequence of physical actions to reach the target location.

This problem is addressed in the field of Vision-and-Language Navigation (VLN) (Anderson et al., 2018b), which focuses on developing autonomous robotic agents that can navigate unseen environments based on natural language instructions. A common VLN benchmark and dataset is Room-to-Room (R2R) (Anderson et al., 2018b), which contains thousands of trajectory–instruction pairs, where the task is to follow natural language instructions to reach a target location. R2R is typically used in combination with the Matterport3D simulator (Anderson et al., 2018b), which simulates indoor environments reconstructed from real-world 3D scans from the Matterport3D dataset (Chang et al., 2018). The simulator represents these environments as navigation graphs, where nodes correspond to navigable locations and edges define transitions between them.

Early approaches to VLN primarily relied on RNN-based sequence-to-sequence models to encode route instructions and predict actions (Anderson et al., 2018b; Fried et al., 2018). Later work shifted toward using pre-trained transformer-based models (Vaswani et al., 2017), which offered improved language understanding and generalization (Li et al., 2019; Chen et al., 2021, 2022).

More recently, researchers have begun exploring the use of Large Language Models (LLMs) and Large Vision-Language Models (LVLMs) for VLN, using both zero-shot prompting (Zhou et al., 2024; Chen et al., 2024) and trained approaches (Zheng et al., 2024; Zhou et al., 2025). While zero-shot methods have shown promise in navigation tasks, their performance still falls short of VLN-specialized transformer-based models (Zhou et al., 2025). Most existing VLN approaches thus seek to train LLMs and LVLMs directly on VLN datasets. Although these trained approaches have achieved

(a) Low-level action space

(b) Panoramic action space

Figure 1: Overview of the approach, which is based on fine-tuning a pre-trained LVLM (Qwen2.5-VL) on the R2R dataset. The model receives as input a multimodal prompt consisting of the route instruction, navigation history, and current view, and outputs the next navigation action.

strong results, they typically rely on custom models that require either changes to the underlying neural architecture or the addition of task-specific components, such as simulators employed at training time (Zheng et al., 2024; Zhou et al., 2025). As a result, the potential of off-the-shelf LVLMs, fine-tuned for VLN without architectural changes, remains largely underexplored.

In addition, the choice of action space – i.e. the possible outputs that the model is designed to generate – has been shown to significantly affect performance (Fried et al., 2018). Early RNN-based approaches typically employed a low-level action space, where the agent observes the environment through an egocentric image and selects from a discrete set of atomic actions such as Move Forward, Turn Left, or Turn Right (Anderson et al., 2018b; Landi et al., 2019). However, low-level action spaces have largely been abandoned in recent work in favor of panoramic action spaces (Li et al., 2019; Chen et al., 2022; Zhou et al., 2024), where the agent perceives its surroundings through a $360°$ panoramic image and chooses among a set of navigable candidate directions, each typically corresponding to an adjacent node in the navigation graph. This shift has been shown to substantially improve performance over low-level alternatives (Fried et al., 2018). While this difference in performance has been explored in the context of RNN-based models (Fried et al., 2018; Landi et al., 2019), it has to our knowledge never been investigated for LVLM-based approaches. While

panoramic action spaces do seem to improve the navigation performance, they also assume a greater prior knowledge about the environment – such as which directions are navigable – and effectively reduce the task to a visually guided graph search (Landi et al., 2019; Krantz et al., 2020). Panoramic action spaces also depend on the availability of panoramic visual input, which in practice requires specialized robot-mounted hardware, such as panoramic or multi-camera rigs.

This paper seeks to address these knowledge gaps through experiments with a state-of-the-art LVLM, Qwen2.5-VL (Bai et al., 2025). An overview of our approach is illustrated in Figure 1. The two main contributions of this paper are:

- The evaluation of off-the-shelf LVLMs (without architectural changes or simulation-based training methods) on VLN through experiments on the R2R dataset.

- An analysis of how the choice of action space (low-level versus panoramic) affects the navigation performance.

The rest of this paper is as follows. We first review related work on Vision-and-Language Navigation and LVLMs. We then present our approach in Section 3, focusing on the fine-tuning process and the definition of possible action spaces. Section 4 then describes the experimental setup and the results obtained on the R2R dataset. Finally, Section 5 discusses those results and Section 6 concludes this paper.

## 2 Related Work

**Large Vision-Language Models in VLN**

Motivated by recent progress with LLMs and LVLMs, several studies have investigated how those models can be applied for VLN. NavGPT (Zhou et al., 2024) employs GPT-4 (OpenAI et al., 2024) in a zero-shot setting, relying on a separate model to convert visual inputs into textual descriptions. In contrast, MapGPT (Chen et al., 2024) prompts GPT-4V to perform joint reasoning over visual inputs and navigation instructions.

NaviLLM (Zheng et al., 2024) uses a frozen Vision Transformer (ViT) (Dosovitskiy et al., 2021) and models spatial relationships between different viewpoints through a trained transformer-based multi-view fusion component which produces a single visual feature for each image. NavGPT-2 (Zhou et al., 2025) uses a frozen LVLM to produce reasoning text from image-instruction pairs and fine-tunes a separate graph-based policy to predict actions and model the topological graph on the fly. Both approaches achieve state-of-the-art performance on R2R, demonstrating the potential of LLMs and LVLMs for navigation.

**Action Spaces in VLN**

Early approaches to VLN employ a low-level action space where the agent perceives the world through an egocentric image at each step and predicts actions such as Move Forward or Turn Right (Anderson et al., 2018b; Wang et al., 2018; Fried et al., 2018). Fried et al., 2018 introduce panoramic action space for VLN. Instead of receiving an egocentric image as input, the model is provided with a panorama comprised of 36 images at different angles. The images closest to the center of an adjacent node are considered as candidate views. Instead of predicting low-level actions, the agent selects between which of these views to navigate to. Using an LSTM (Hochreiter and Schmidhuber, 1997) seq-2-seq model, they observe a 12% performance increase on R2R when going from low-level to panoramic action space.

Although there is little recent work on low-level action spaces in discrete environments (VLN-DE), it remains the most common approach for VLN in continuous environments (VLN-CE) (Krantz et al., 2020; Zhang et al., 2024) where agents are tasked with navigating environments not constrained by a predefined navigation graph. In this work, we focus on VLN in discrete environments.

**Modality alignment in LVLMs**

Modern Large Vision-Language Models (LVLMs) typically comprise three core components: a vision encoder (e.g., a Vision Transformer (Dosovitskiy et al., 2021)), a cross-modal projector, and a text encoder (e.g, an LLM) (Bai et al., 2025). The role of the cross-modal projector is to align the visual features produced by the vision encoder with the latent space of the LLM.

Laurençon et al. (2024) investigate key design choices in building LVLMs and identify two prevalent architectural paradigms for vision-language alignment. The first is the *cross-attention architectures*, in which visual features are injected at different layers within the LLM, one example of such a model is Flamingo (Alayrac et al., 2022). The second is the *fully autoregressive architectures* where the output of the vision encoder is projected into the input space of the LLM and concatenated with the sequence of text embeddings as a multimodal prompt (Zhu et al., 2023; Li et al., 2023). The model used in this study, Qwen2.5-VL, follows this fully autoregressive design.

## 3 Method

### 3.1 Problem Formulation

We adopt the standard VLN in discrete environments (VLN-DE) setup (Anderson et al., 2018b; Fried et al., 2018; Chen et al., 2022), where the environment is modeled as an undirected graph $G = \{V, E\}$. The nodes $V = \{v_i\}_{i=1}^K$ represent $K$ navigable locations while the edges $E$ constitute navigation paths between them. We then formulate the problem of following route instructions in a graph-based environment as follows: given a natural language route instruction $W = \{w_1, w_2, \ldots, w_L\}$, the agent is tasked with following the instruction to reach the goal location. At each time step $t$, the agent receives a visual observation $O_t$, maintains a history context $H_t$, and is provided with auxiliary signals such as the cumulative distance traveled $d_t \in \mathbb{R}$ and the current step number $t$. The specific formulation of the agent's input and output depends on the underlying action space, as described below.

### 3.2 Low-level Action Space

In the low-level action space, the agent perceives its environment through an egocentric image $O_t$ at each step. It maintains a historical context $H_t = \{(O_1, a_1), (O_2, a_2), \ldots, (O_{t-1}, a_{t-1})\}$

where $O_{t-1}$ and $a_{t-1}$ are the image and action from the previous step, respectively. Additionally, the agent is provided with a set of low-level actions $U_t = \{u_1, u_2, \ldots, u_k\}$ that represent the actions allowed at step $t$, given the physical constraints of the environment (e.g., the agent cannot move forward if directly facing a wall). The agent predicts the next action $a_t$ by estimating the probability:

$$P(a_t \mid W, O_t, H_t, d_t, t, U_t) \qquad (1)$$

The low-level action space used in this work consists of four discrete actions:

- Move: moves forward to the node closest to the center of the current field of view.

- Left, Right: rotate the agent by $30°$ in the respective direction.

- Stop: signals that the agent believes it has reached the goal.

A limitation of this setup is that navigation is constrained to a discrete graph of nodes. The Move action advances the agent to the node most centered in its current field of view, but this target is not necessarily aligned with the agent's heading. As a result, the agent may appear to move sideways, which can lead to non-intuitive trajectories. To mitigate this, an automatic reorientation step, referred to as Automatically Turn Towards Node, is applied before each Move action. Although this reorientation is not part of the learnable action space, both the resulting observation and action are included in the agent's history. This adjustment allows us to evaluate whether explicitly aligning the agent's heading with its movement direction improves navigation performance.

### 3.3 Panoramic Action Space

With the panoramic action space, the agent perceives the environment through a $360°$ panoramic image $O_t$ at each step, aligned with its current heading. The agent maintains a history of panoramic views $H_t = \{O_1, \ldots, O_{t-1}\}$ and selects from a set of navigable candidate views $C_t = \{c_1, \ldots, c_k\}$. Each candidate $c_i$ includes an image, a relative heading $\theta_i \in [-180°, 180°]$, and an associated travel distance $\delta_i \in \mathbb{R}_{\geq 0}$. The task for each step is to predict the correct candidate direction $c_t$:

$$P(c_t \mid W, O_t, H_t, d_t, t, C_t) \qquad (2)$$

Similarly to low-level actions, the episode concludes when the agent predicts the Stop action.

The panoramic image is centered on the agent's current heading, while each candidate view is a standard perspective image oriented directly toward a navigable direction. Candidate views are sorted from left to right based on their relative angle to the panoramic center, with the leftmost candidate assigned index 0 and the rightmost index $K - 1$.

At each step, the model predicts a token corresponding to one of the candidate indices (from 0 to $K-1$) or the Stop action. Unlike traditional panoramic setups (Fried et al., 2018; Zheng et al., 2024), where candidate views are extracted from within the panorama itself, this approach treats the panorama and candidate views as separate inputs. This design, motivated by memory limitations, reduces the number of input images per step. See Appendix A for an illustrative example.

### 3.4 Action selection

To select the next action to perform, the model receives a structured multi-modal prompt that encodes the current state, including the instruction, visual input, and auxiliary information such as step number and distance traveled. These prompts follow a fixed schema shown in Figure 2. Inference is performed greedily, selecting the most probable action at each step without backtracking.

In addition to the dynamic input state, each prompt includes a static system prompt that explains the task and describes the individual input fields. The system prompt is fixed and specific to each action space, and remains unchanged throughout training and evaluation. The full system prompts are included in Appendix A

### 3.5 Fine-tuning

The LVLMs are fine-tuned through behavior cloning, where the model learns to imitate expert demonstrations. At each time step $t$, the model receives a multimodal prompt $x_t$ represents the current state, and is trained to predict the expert action $a_t$ as a token from its own vocabulary. The training objective minimizes the total negative log-likelihood of the expert actions over the entire episode. Gradients are accumulated across all time steps in an episode, and the weights are updated at the end of each episode.

Unlike many recent VLN approaches (Chen et al., 2021; Zhou et al., 2025; Anderson et al.,

(a) Low-level action space prompt schema      (b) Panoramic action space prompt schema

Figure 2: Prompt schemata for low-level and panoramic action spaces.

2018b), our approach does not therefore rely on reinforcement learning or student forcing, but simply fine-tunes the LVLM model on the basis of expert routes. A key advantage of this approach is the fact that it can be applied without access to a simulator at training time.

## 4 Experiments

The proposed approach was evaluated on the Room-to-Room (R2R) dataset using both *offline* and *online* evaluation modes. The offline mode assesses the model's ability to follow expert trajectories, whereas the online mode evaluates its performance when navigating autonomously within the Matterport3D simulation environment.

### 4.1 Dataset

The Room-to-Room dataset (Anderson et al., 2018b) contains 21,567 English route instructions corresponding to 7,189 trajectories across 90 environments. Each ground truth trajectory is a sequence of nodes in a Matterport3D environment. Each trajectory has 3 corresponding instructions.

The dataset is split into four subsets: *train* (61 environments), *val seen* (56 environments overlapping with train), *val unseen* (11 environments), and *test* (18 environments). Performance is evaluated on the val unseen and test splits. All splits are preprocessed to convert ground truth trajectories into sequences of actions. [1]

---

### 4.2 Evaluation Metrics

Online, the models are evaluated using standard VLN metrics (Anderson et al., 2018b). **Navigation Error (NE)** is the average walkable distance between the agent's final location and the goal location in meters. An episode is considered successful if NE $\leq 3$ m and the last predicted action is Stop. **Path Length (PL)** is the average path length (in meters). **Oracle Success Rate (OSR)** measures the percentage of episodes in which the agent was within 3 meters of the goal at any point during the navigation episode. **Success Rate (SR)** is the percentage of episodes that are successful. **Success Weighted by Path Length (SPL)** (Anderson et al., 2018a) combines SR with path length, penalizing unnecessarily long paths. **Coverage Weighted by Length Score (CLS)** (Jain et al., 2019) measures how well the agent's predicted path follows the route instruction, penalizing paths that deviate from the ground truth path.

For the offline evaluation, the reported metrics are **Accuracy** the proportion of actions correctly predicted by the model; **Macro F1**, the unweighted mean F1 score computed across all action classes; and **Conservative Success Rate (CSR)**, the percentage of episodes in which all actions are identical (from start to finish) to the actions selected by the expert. For offline evaluation, we use the third instruction from each trajectory.

---

### 4.3 Implementation

#### Model

We experimented with two distinct Large Vision-Language Models (LVLMs): Qwen2-VL-2B-Instruct (Wang et al., 2024) and Qwen2.5-VL-3B-Instruct (Bai et al., 2025). Qwen2.5-VL-3B is the larger of the two and is pre-trained on 4 trillion tokens, compared to 1.2 trillion tokens for Qwen2-VL-2B. As our experiments showed that Qwen2.5 consistently provided superior performance compared to Qwen2 in both offline and online metrics on the validation dataset, we only provide evaluation results obtained for Qwen2.5.

During fine-tuning, the vision encoder and the cross-modal projection layer are kept frozen, as preliminary experiments indicated that tuning only the LLM led to improved performance.[2]

#### Simulator

The Matterport3D simulator (MP3D) is used for evaluation and for generating the preprocessed training data. In MP3D, the agent's field of view is determined by the image resolution and the vertical field of view (VFOV). This work uses an image resolution of $640 \times 480$ for egocentric and candidate images. The VFOV is set to $105°$ to allow the agent to perceive a broader visual context. This is substantially larger than the VFOV used in prior work, which typically ranges from $60°$ (Anderson et al., 2018b; Fried et al., 2018) to $75°$ (Krantz et al., 2020). Panoramic images are constructed by stitching together three egocentric views captured while rotating the agent in place.

#### Training

All models are fine-tuned with a batch size of 1, a learning rate of 1e-5, and a weight decay of 0.1. A linear learning rate scheduler is used with warmup over the first 10% of training steps. FlashAttention (Dao et al., 2022) is enabled, and training is performed in bfloat16 precision. Input images are resized to half their original size to accommodate GPU memory constraints.[3] Experiments were conducted on a single NVIDIA A100 80GB GPU. Models are fine-tuned for 1 epoch across all instructions, corresponding to 3 total passes over the unique paths (as each path in R2R is associated

---

[2]The trained models are publicly available at `https://huggingface.co/Vebbern` for reproducibility.

[3]meaning $320 \times 240$ for candidates and egocentric views, and $960 \times 240$ for panoramic views

| Model | Accuracy↑ | Macro F1↑ | CSR↑ |
|---|---|---|---|
| **Val seen:** | | | |
| Qwen2.5-VL-low | 0.73 | 0.74 | 0.04 |
| Qwen2.5-VL-pano | 0.73 | 0.61 | 0.16 |
| **Val unseen:** | | | |
| Qwen2.5-VL-low | 0.73 | 0.73 | 0.03 |
| Qwen2.5-VL-pano | 0.73 | 0.62 | 0.15 |

Table 1: Offline evaluation results on the seen and unseen R2R validation sets.

with three distinct route instructions).[4]

### 4.4 Results

**Offline evaluation** Table 1 presents the offline evaluation results after fine-tuning the Qwen2.5 model on the full training set of R2R. In terms of accuracy, the panoramic and low-level models score similarly. The low-level model has a higher macro F1 score, which could be explained by the larger number of actions of panoramic models (up to 12 candidate views). However, the panoramic model has a significantly higher conservative success rate (CSR) than the low-level one. Qwen2.5-VL-pano achieves a CSR of 15% on val unseen, compared to a mere 3% CSR for Qwen2.5-VL-low.

**Online evaluation** Table 2 compares our results with state-of-the-art (SOTA) approaches on R2R using single-run greedy search (i.e., no pre-exploration). Results are shown for both panoramic and low-level action space.

The model fine-tuned for low-level action spaces, Qwen2.5-VL-low, achieves a success rate (SR) of 26% on the test set, outperforming the original R2R baseline (Seq2Seq, 21% SR) and Speaker-Follower (SF) without panoramic action (25% SR on val unseen). However, it still lags behind the LSTM-based DCF model of (Landi et al., 2019), which reached 35% SR, despite being substantially smaller in size. However, Qwen2.5-VL-low is less prone to overfitting to training environments, as reflected in the smaller SR gap between val seen and unseen (35% vs. 27%) compared to DCF, which drops from 58% to 34% on val unseen.

The panoramic model, Qwen2.5-VL-pano, achieves a 41% SR on the R2R test set. This outperforms all low-level models as well as panoramic approaches such as Speaker-Follower (36% on val unseen) and NavGPT (Zhou et al., 2024) (34% on

---

[4]The source code is available at `https://github.com/Vebjorhk/masters-thesis-VLN`.

| | **Val Seen** | | | | | **Val Unseen** | | | | | **Test (Unseen)** | | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | PL | NE↓ | OSR↑ | SR↑ | SPL↑ | PL | NE↓ | OSR↑ | SR↑ | SPL↑ | PL | NE↓ | OSR↑ | SR↑ | SPL↑ |
| Human | - | - | - | - | - | - | - | - | - | - | 11.85 | 1.61 | 90 | 86 | 76 |
| **Low-Level** | | | | | | | | | | | | | | | |
| Seq2Seq (2018b) | 11.33 | 6.01 | 53 | 39 | - | 8.39 | 7.81 | 28 | 22 | - | 8.13 | 7.85 | 27 | 21 | - |
| SF (2018) | - | 4.28 | 60 | 47 | - | - | 5.75 | 33 | 25 | - | - | - | - | - | - |
| RPA(2018) | 8.46 | 5.56 | 53 | 43 | - | 7.22 | 7.65 | 32 | 25 | - | 9.15 | 7.53 | 33 | 25 | - |
| DCF (2019) | - | 3.96 | 73 | 58 | 51 | - | 6.52 | 43 | 34 | 29 | 9.81 | 6.55 | 45 | 35 | 31 |
| **Panoramic** | | | | | | | | | | | | | | | |
| SF (2018) | - | 3.36 | 73 | 66 | - | - | 6.62 | 45 | 36 | - | - | - | - | - | - |
| PRESS (2019) | 10.35 | 3.09 | - | 71 | 67 | 10.06 | 4.31 | - | 59 | 55 | 10.52 | 4.53 | 63 | 57 | 53 |
| VLN ↻ BERT (2021) | 11.13 | 2.90 | - | 72 | 68 | 12.01 | 3.93 | - | 63 | 57 | 12.35 | 4.09 | - | 63 | 57 |
| HAMT (2021) | 11.15 | 2.51 | - | 76 | 72 | 11.46 | 2.29 | - | 66 | 61 | 12.27 | 3.93 | - | 65 | 60 |
| DUET (2022) | - | - | - | - | - | 13.94 | 3.31 | - | 72 | 60 | 14.73 | 3.65 | - | 69 | 59 |
| NavGPT (2024) | - | - | - | - | - | - | - | - | - | - | 11.45 | 6.46 | 42 | 34 | 29 |
| NaviLLM (2024) | - | - | - | - | - | - | - | - | - | 59 | - | - | - | - | 60 |
| NavGPT-2 (2025) | 14.13 | 2.84 | 83 | 74 | 63 | 14.01 | 2.98 | 84 | 74 | 61 | 14.74 | 3.33 | 80 | 72 | 60 |
| **Qwen2.5-VL-low** | 10.27 | 7.14 | 41 | 35 | 32 | 10.50 | 7.84 | 34 | 27 | 24 | 10.59 | 7.99 | 34 | 26 | 24 |
| **Qwen2.5-VL-pano** | 9.98 | 5.69 | 56 | 50 | 47 | 9.83 | 6.65 | 46 | 38 | 35 | 9.96 | 6.53 | 50 | 41 | 38 |

Table 2: Comparison of panoramic and low-level models with state-of-the-art performance using single-run greedy search. CLS is not reported for R2R test set. Models introduced in this work are shown in bold.

| Models | PL | NE↓ | OSR↑ | SR↑ | SPL↑ | CLS↑ |
|---|---|---|---|---|---|---|
| **Val Unseen:** | | | | | | |
| 105-VFOV | 10.17 | 7.87 | 0.34 | 0.25 | 0.23 | 0.45 |
| 82-VFOV | 9.9 | 7.87 | 0.32 | 0.25 | 0.23 | 0.44 |
| **No-Adjust** | 10.72 | 7.7 | 0.38 | 0.29 | 0.26 | 0.44 |

Table 3: Online results on R2R val unseen for alternative definitions of the low-level action space.

test). However, this model falls short of more recent task-specific panoramic approaches such as NaviLLM (Zheng et al., 2024) (60% SPL) and NavGPT-2 (Zhou et al., 2025) (72% SR).

**Alternative low-level action spaces** We also explored alternative configurations for the low-level action space. Specifically, we assessed the impact of (1) disabling the `Automatically Turn Towards Node` action, and (2) reducing the vertical field of view (VFOV) from $105°$ to a narrower $82°$.

Table 3 presents the performance on the val unseen split for those two alternatives. The difference between $82°$ and $105°$ VFOV is minimal, with only slight improvements in CLS and OSR scores for the $105°$ configuration. However, removing the adjustment action leads to a noticeable performance gain: the No-Adjust model achieves a SR of 29%, compared to 25% for the default. This suggests that explicitly facing the next node before movement is often unnecessary for effective navigation.

## 5 Discussion

**Fine-tuning off-the-shelf LVLMs for R2R** The results indicate that fine-tuning off-the-shelf LVLMs such as Qwen2.5-VL on the R2R task fails to yield strong performance, despite the fact that such models are significantly larger than older, VLN-specific architectures such as PRESS (Li et al., 2019), DUET (Chen et al., 2022), and HAMT (Chen et al., 2021). It is difficult to pinpoint the exact source of this performance gap, though our use of behavior cloning – rather than optimization through student forcing and/or reinforcement learning, as done by e.g. (Chen et al., 2021; Zhou et al., 2025) – may be a contributing factor.

Compared to NaviLLM (Zheng et al., 2024) and NavGPT-2 (Zhou et al., 2025), which are the two approaches most similar to this work, a key difference becomes apparent. In the Qwen2.5-VL-low model, each input image is encoded and then fed directly into the LLM, which is solely responsible for interpreting the route instruction, modeling spatial relationships between images, and predicting actions. While Qwen2.5-VL reduces visual token count through patch merging, it does not incorporate any explicit mechanisms for modeling spatial structure between images before they are fed to the LLM. In contrast, NaviLLM (Zheng et al., 2024) includes a transformer-based module that explicitly captures the spatial relationships be-

| Split | Avg. steps per path (low-level) | Avg. steps per path (panoramic) |
|---|---|---|
| train | 12.88 | 6.00 |
| val seen | 12.85 | 6.07 |
| val unseen | 13.40 | 5.97 |

Table 4: Average number of steps (actions) for panoramic and low-level variants of R2R.



Figure 3: Avg. path length (meters) on R2R val unseen.

tween panoramic images before it is fed as input to the LLM. NavGPT-2 (Zhou et al., 2025) takes this further by using a separate graph-based policy network to model viewpoint connectivity and predict actions, while delegating route-level reasoning to the LLM. These design differences may help explain at least part of the observed performance gap: relying solely on the LLM for spatial reasoning and control can be challenging – especially for longer paths – compared to models that explicitly encode spatial structure. Prior work also suggests that reducing visual tokens benefits non-OCR tasks (Laurençon et al., 2024). Both Nav-iLLM and NavGPT-2 use significantly fewer visual tokens than Qwen2.5-VL (Bai et al., 2025).

**Panoramic vs. low-level action space** The panoramic models consistently outperform low-level ones, which aligns with previous findings by Fried et al. (2018), although the performance gap in our setup is slightly larger (16% vs. 12% SR) showing that the panoramic approach leads to better results for LVLMs as well. One plausible explanation for the performance gap is that low-level action sequences are, on average, twice as long as those in the panoramic setting (Table 4). As shown in Figure 3, both model types tend to perform better on shorter trajectories. This suggests that longer sequences in the low-level setting increase the diffi-

culty of the task, as they provide more opportunities for errors to accumulate and make recovery more challenging. This is further supported by the noticeably higher Conservative Success Rate (CSR) for panoramic models (Table 1), indicating they are more likely to keep the agent on the correct path. In contrast, low-level models are more prone to errors due to the increased number of decision points, making it harder to recover once the agent deviates from the intended path.

The extent to which the additional visual information provided by panoramic images contributes to improved performance remains somewhat unclear. Panoramic observations may benefit the agent by reducing the need for physical reorientation to perceive important landmarks. Low-level action spaces may also place greater demands on spatial and temporal reasoning abilities: the agent must not only ground instructions in the visual context but also anticipate when certain actions should be executed – such as recognizing that a given action may only occur after completing several turns.

## 6 Conclusion

This work focused on the use of off-the-shelf Large Vision-Language Models (LVLMs) for Vision-and-Language Navigation (VLN) tasks. More precisely, we investigated how such models could be fine-tuned directly from expert routes, without modifying the model's underlying architecture or relying on online approaches that necessitate the use of a simulator at training time. The performance of this approach was assessed through experiments on the R2R dataset and explored using both panoramic and low-level action spaces.

The best performing model, fine-tuned from Qwen2.5-VL, achieved a success rate (SR) of 41% on the R2R dataset. Our results suggest that simply fine-tuning LVLMs remains insufficient to reach state-of-the-art performance on navigation tasks. Furthermore, we find that the performance gap between low-level and panoramic action spaces persists even with larger and more powerful models, with a 16% difference in SR on the R2R test set in favor of the panoramic setup.

A promising topic for future work is the systematic study of off-the-shelf LVLMs on the R2R dataset. Evaluating a broader range of models beyond Qwen2 and Qwen2.5 could help identify which architectural choices lead to better navigation performance. Additionally, a more focused

investigation of the panoramic action space is warranted – particularly through ablation studies that isolate the effect of including a panoramic view, and systematically vary the field of view to understand its impact on performance. We also encourage future work to further explore the low-level action space for more recent approaches, including adapting it to existing state-of-the-art methods such as NaviLLM (Zheng et al., 2024) and NavGPT-2 (Zhou et al., 2025) and comparing the performance to panoramic action space.

## Limitations

We acknowledge several limitations in this work. Most importantly, the fine-tuning approach is limited to behavior cloning, and did not include the use of VLN training techniques such as student forcing or reinforcement learning, potentially limiting direct comparability with prior work that leverages these strategies. For evaluation, we set up a web API to communicate with the machine running the simulator remotely. However, this introduced an additional limitation: the need for network calls made simulator evaluation significantly more time-consuming. As a result, we restricted evaluation of alternative setups to only the first third of the route instructions.

GPU memory limitations restricted training to a batch size of 1. To further reduce memory usage, we deviated from the standard panoramic action format used in many VLN approaches (Fried et al., 2018; Li et al., 2019; Hong et al., 2021), where the model receives a set of discrete view images (typically 36) and selects candidate views from among them. Instead, we preprocessed the full panorama as a single image and treated candidate views as separate, independent inputs. This setup reduces granularity, introduces visual artifacts, and limits comparability to prior benchmarks.

Finally, we note that Room-to-Room (R2R) contains only English-language route instructions, which limits the applicability of our approach to English-only scenarios. While multilingual VLN datasets have been proposed – such as Room-across-Room (RxR) (Ku et al., 2020) – our current experiments do not address multilingual aspects.

## Ethics Statement

This work investigated the use of off-the-shelf Large Vision-Language Models (LVLMs) for Vision-and-Language Navigation (VLN) tasks. All models used in this study are open-source and publicly available. The dataset employed, Room-to-Room (R2R), is a widely used benchmark in the VLN community and does not contain personally identifiable information. We do not foresee any direct ethical concerns related to the methods, data, or potential applications of this research. Our study adheres to the ACL Code of Ethics.

## References

Jean-Baptiste Alayrac, Jeff Donahue, Pauline Luc, Antoine Miech, Iain Barr, Yana Hasson, Karel Lenc, Arthur Mensch, Katherine Millican, Malcolm Reynolds, Roman Ring, Eliza Rutherford, Serkan Cabi, Tengda Han, Zhitao Gong, Sina Samangooei, Marianne Monteiro, Jacob L. Menick, Sebastian Borgeaud, Andy Brock, Aida Nematzadeh, Sahand Sharifzadeh, Mikołaj Bińkowski, Ricardo Barreira, Oriol Vinyals, Andrew Zisserman, and Karén Simonyan. 2022. Flamingo: a Visual Language Model for Few-Shot Learning. In *Advances in Neural Information Processing Systems*, volume 35, pages 23716–23736.

Peter Anderson, Angel Chang, Devendra Singh Chaplot, Alexey Dosovitskiy, Saurabh Gupta, Vladlen Koltun, Jana Kosecka, Jitendra Malik, Roozbeh Mottaghi, Manolis Savva, and Amir R. Zamir. 2018a. On Evaluation of Embodied Navigation Agents. ArXiv:1807.06757.

Peter Anderson, Qi Wu, Damien Teney, Jake Bruce, Mark Johnson, Niko Sunderhauf, Ian Reid, Stephen Gould, and Anton van den Hengel. 2018b. Vision-and-Language Navigation: Interpreting Visually-Grounded Navigation Instructions in Real Environments. In *2018 IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 3674–3683.

Shuai Bai, Keqin Chen, Xuejing Liu, Jialin Wang, Wenbin Ge, Sibo Song, Kai Dang, Peng Wang, Shijie Wang, Jun Tang, Humen Zhong, Yuanzhi Zhu, Mingkun Yang, Zhaohai Li, Jianqiang Wan, Pengfei Wang, Wei Ding, Zheren Fu, Yiheng Xu, Jiabo Ye, Xi Zhang, Tianbao Xie, Zesen Cheng, Hang Zhang, Zhibo Yang, Haiyang Xu, and Junyang Lin. 2025. Qwen2.5-VL Technical Report. ArXiv:2502.13923.

Angel Chang, Angela Dai, Thomas Funkhouser, Maciej Halber, Matthias Niebner, Manolis Savva, Shuran Song, Andy Zeng, and Yinda Zhang. 2018. Matterport3D: Learning from RGB-D data in indoor environments. In *Proceedings - 2017 International Conference on 3D Vision, 3DV 2017*, pages 667–676.

Jiaqi Chen, Bingqian Lin, Ran Xu, Zhenhua Chai, Xiaodan Liang, and Kwan-Yee Wong. 2024. MapGPT: Map-Guided Prompting with Adaptive Path Planning for Vision-and-Language Navigation. In *Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 9796–9810, Bangkok, Thailand. Association for Computational Linguistics.

Shizhe Chen, Pierre-Louis Guhur, Cordelia Schmid, and Ivan Laptev. 2021. History Aware Multimodal Transformer for Vision-and-Language Navigation. In *Advances in Neural Information Processing Systems*, volume 34, pages 5834–5847. Curran Associates, Inc.

Shizhe Chen, Pierre-Louis Guhur, Makarand Tapaswi, Cordelia Schmid, and Ivan Laptev. 2022. Think Global, Act Local: Dual-Scale Graph Transformer for Vision-and-Language Navigation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 16537–16547.

Tri Dao, Dan Fu, Stefano Ermon, Atri Rudra, and Christopher Ré. 2022. FlashAttention: Fast and Memory-Efficient Exact Attention with IO-Awareness. In *Advances in Neural Information Processing Systems*, volume 35, pages 16344–16359.

Alexey Dosovitskiy, Lucas Beyer, Alexander Kolesnikov, Dirk Weissenborn, Xiaohua Zhai, Thomas Unterthiner, Mostafa Dehghani, Matthias Minderer, Georg Heigold, Sylvain Gelly, Jakob Uszkoreit, and Neil Houlsby. 2021. An Image is Worth 16x16 Words: Transformers for Image Recognition at Scale. In *International Conference on Learning Representations*.

Daniel Fried, Ronghang Hu, Volkan Cirik, Anna Rohrbach, Jacob Andreas, Louis-Philippe Morency, Taylor Berg-Kirkpatrick, Kate Saenko, Dan Klein, and Trevor Darrell. 2018. Speaker-Follower Models for Vision-and-Language Navigation. In *Advances in Neural Information Processing Systems*, volume 31. Curran Associates, Inc.

Sepp Hochreiter and Jürgen Schmidhuber. 1997. Long Short-term Memory. *Neural computation*, 9:1735–80.

Yicong Hong, Qi Wu, Yuankai Qi, Cristian Rodriguez-Opazo, and Stephen Gould. 2021. VLNBERT: A Recurrent Vision-and-Language BERT for Navigation. In *2021 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 1643–1653.

Vihan Jain, Gabriel Magalhaes, Alexander Ku, Ashish Vaswani, Eugene Ie, and Jason Baldridge. 2019. Stay on the Path: Instruction Fidelity in Vision-and-Language Navigation. In *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*, pages 1862–1872, Florence, Italy. Association for Computational Linguistics.

Jacob Krantz, Erik Wijmans, Arjun Majumdar, Dhruv Batra, and Stefan Lee. 2020. Beyond the Nav-Graph: Vision-and-Language Navigation in Continuous Environments. In *Computer Vision – ECCV 2020*, pages 104–120, Cham. Springer International Publishing.

Alexander Ku, Peter Anderson, Roma Patel, Eugene Ie, and Jason Baldridge. 2020. Room-Across-Room: Multilingual Vision-and-Language Navigation with Dense Spatiotemporal Grounding. In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pages 4392–4412, Online. Association for Computational Linguistics.

Federico Landi, Lorenzo Baraldi, Massimiliano Corsini, and Rita Cucchiara. 2019. Embodied Vision-and-Language Navigation with Dynamic Convolutional Filters. In *Proceedings of the British Machine Vision Conference*.

Hugo Laurençon, Léo Tronchon, Matthieu Cord, and Victor Sanh. 2024. What matters when building vision-language models? In *Advances in Neural Information Processing Systems*, volume 37, pages 87874–87907.

Junnan Li, Dongxu Li, Silvio Savarese, and Steven Hoi. 2023. BLIP-2: Bootstrapping Language-Image Pre-training with Frozen Image Encoders and Large Language Models. In *Proceedings of the 40th International Conference on Machine Learning*, pages 19730–19742.

Xiujun Li, Chunyuan Li, Qiaolin Xia, Yonatan Bisk, Asli Celikyilmaz, Jianfeng Gao, Noah A. Smith, and Yejin Choi. 2019. Robust Navigation with Language Pretraining and Stochastic Sampling. In *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)*, pages 1494–1499, Hong Kong, China. Association for Computational Linguistics.

OpenAI, Josh Achiam, Steven Adler, Sandhini Agarwal, Lama Ahmad, Ilge Akkaya, Florencia Leoni Aleman, Diogo Almeida, Janko Altenschmidt, Sam Altman, Shyamal Anadkat, Red Avila, Igor Babuschkin, Suchir Balaji, Valerie Balcom, Paul Baltescu, Haiming Bao, Mohammad Bavarian, Jeff Belgum, Irwan Bello, Jake Berdine, Gabriel Bernadett-Shapiro, Christopher Berner, Lenny Bogdonoff, Oleg Boiko, Madelaine Boyd, Anna-Luisa Brakman, Greg Brockman, Tim Brooks, Miles Brundage, Kevin Button, Trevor Cai, Rosie Campbell, Andrew Cann, Brittany Carey, Chelsea Carlson, Rory Carmichael, Brooke Chan, Che Chang, Fotis Chantzis, Derek Chen, Sully

458

Chen, Ruby Chen, Jason Chen, Mark Chen, Ben Chess, Chester Cho, Casey Chu, Hyung Won Chung, Dave Cummings, Jeremiah Currier, Yunxing Dai, Cory Decareaux, Thomas Degry, Noah Deutsch, Damien Deville, Arka Dhar, David Dohan, Steve Dowling, Sheila Dunning, Adrien Ecoffet, Atty Eleti, Tyna Eloundou, David Farhi, Liam Fedus, Niko Felix, Simón Posada Fishman, Juston Forte, Isabella Fulford, Leo Gao, Elie Georges, Christian Gibson, Vik Goel, Tarun Gogineni, Gabriel Goh, Rapha Gontijo-Lopes, Jonathan Gordon, Morgan Grafstein, Scott Gray, Ryan Greene, Joshua Gross, Shixiang Shane Gu, Yufei Guo, Chris Hallacy, Jesse Han, Jeff Harris, Yuchen He, Mike Heaton, Johannes Heidecke, Chris Hesse, Alan Hickey, Wade Hickey, Peter Hoeschele, Brandon Houghton, Kenny Hsu, Shengli Hu, Xin Hu, Joost Huizinga, Shantanu Jain, Shawn Jain, Joanne Jang, Angela Jiang, Roger Jiang, Haozhun Jin, Denny Jin, Shino Jomoto, Billie Jonn, Heewoo Jun, Tomer Kaftan, Łukasz Kaiser, Ali Kamali, Ingmar Kanitscheider, Nitish Shirish Keskar, Tabarak Khan, Logan Kilpatrick, Jong Wook Kim, Christina Kim, Yongjik Kim, Jan Hendrik Kirchner, Jamie Kiros, Matt Knight, Daniel Kokotajlo, Łukasz Kondraciuk, Andrew Kondrich, Aris Konstantinidis, Kyle Kosic, Gretchen Krueger, Vishal Kuo, Michael Lampe, Ikai Lan, Teddy Lee, Jan Leike, Jade Leung, Daniel Levy, Chak Ming Li, Rachel Lim, Molly Lin, Stephanie Lin, Mateusz Litwin, Theresa Lopez, Ryan Lowe, Patricia Lue, Anna Makanju, Kim Malfacini, Sam Manning, Todor Markov, Yaniv Markovski, Bianca Martin, Katie Mayer, Andrew Mayne, Bob McGrew, Scott Mayer McKinney, Christine McLeavey, Paul McMillan, Jake McNeil, David Medina, Aalok Mehta, Jacob Menick, Luke Metz, Andrey Mishchenko, Pamela Mishkin, Vinnie Monaco, Evan Morikawa, Daniel Mossing, Tong Mu, Mira Murati, Oleg Murk, David Mély, Ashvin Nair, Reiichiro Nakano, Rajeev Nayak, Arvind Neelakantan, Richard Ngo, Hyeonwoo Noh, Long Ouyang, Cullen O'Keefe, Jakub Pachocki, Alex Paino, Joe Palermo, Ashley Pantuliano, Giambattista Parascandolo, Joel Parish, Emy Parparita, Alex Passos, Mikhail Pavlov, Andrew Peng, Adam Perelman, Filipe de Avila Belbute Peres, Michael Petrov, Henrique Ponde de Oliveira Pinto, Michael, Pokorny, Michelle Pokrass, Vitchyr H. Pong, Tolly Powell, Alethea Power, Boris Power, Elizabeth Proehl, Raul Puri, Alec Radford, Jack Rae, Aditya Ramesh, Cameron Raymond, Francis Real, Kendra Rimbach, Carl Ross, Bob Rotsted, Henri Roussez, Nick Ryder, Mario Saltarelli, Ted Sanders, Shibani Santurkar, Girish Sastry, Heather Schmidt, David Schnurr, John Schulman, Daniel Selsam, Kyla Sheppard, Toki Sherbakov, Jessica Shieh, Sarah Shoker, Pranav Shyam, Szymon Sidor, Eric Sigler, Maddie Simens, Jordan Sitkin, Katarina Slama, Ian Sohl, Benjamin Sokolowsky, Yang Song, Natalie Staudacher, Felipe Petroski Such, Natalie Summers, Ilya Sutskever, Jie Tang, Nikolas Tezak, Madeleine B. Thompson, Phil Tillet, Amin Tootoonchian, Elizabeth Tseng, Preston Tuggle, Nick Turley, Jerry Tworek, Juan Felipe Cerón Uribe, Andrea Vallone, Arun Vijayvergiya, Chelsea Voss, Carroll Wainwright, Justin Jay Wang, Alvin Wang, Ben Wang, Jonathan Ward, Jason Wei, C. J. Weinmann, Akila Welihinda, Peter Welinder, Jiayi Weng, Lilian Weng, Matt Wiethoff, Dave Willner, Clemens Winter, Samuel Wolrich, Hannah Wong, Lauren Workman, Sherwin Wu, Jeff Wu, Michael Wu, Kai Xiao, Tao Xu, Sarah Yoo, Kevin Yu, Qiming Yuan, Wojciech Zaremba, Rowan Zellers, Chong Zhang, Marvin Zhang, Shengjia Zhao, Tianhao Zheng, Juntang Zhuang, William Zhuk, and Barret Zoph. 2024. GPT-4 Technical Report. ArXiv:2303.08774.

Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Ł ukasz Kaiser, and Illia Polosukhin. 2017. Attention is All you Need. In *Advances in Neural Information Processing Systems*, volume 30. Curran Associates, Inc.

Peng Wang, Shuai Bai, Sinan Tan, Shijie Wang, Zhihao Fan, Jinze Bai, Keqin Chen, Xuejing Liu, Jialin Wang, Wenbin Ge, Yang Fan, Kai Dang, Mengfei Du, Xuancheng Ren, Rui Men, Dayiheng Liu, Chang Zhou, Jingren Zhou, and Junyang Lin. 2024. Qwen2-VL: Enhancing Vision-Language Model's Perception of the World at Any Resolution. ArXiv:2409.12191.

Xin Wang, Wenhan Xiong, Hongmin Wang, and William Yang Wang. 2018. Look Before You Leap: Bridging Model-Free and Model-Based Reinforcement Learning for Planned-Ahead Vision-and-Language Navigation. In *Proceedings of the European Conference on Computer Vision (ECCV)*, pages 37–53.

Jiazhao Zhang, Kunyu Wang, Rongtao Xu, Gengze Zhou, Yicong Hong, Xiaomeng Fang, Qi Wu, Zhizheng Zhang, and He Wang. 2024. NaVid: Video-based VLM Plans the Next Step for Vision-and-Language Navigation. In *Robotics: Science and Systems*.

Duo Zheng, Shijia Huang, Lin Zhao, Yiwu Zhong, and Liwei Wang. 2024. Towards Learning a Generalist Model for Embodied Navigation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 13624–13634.

Gengze Zhou, Yicong Hong, Zun Wang, Xin Eric Wang, and Qi Wu. 2025. NavGPT-2: Unleashing Navigational Reasoning Capability for Large Vision-Language Models. In *Computer Vision – ECCV 2024*, pages 260–278, Cham. Springer Nature Switzerland.

Gengze Zhou, Yicong Hong, and Qi Wu. 2024. NavGPT: Explicit Reasoning in Vision-and-Language Navigation with Large Language Models. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 38, pages 7641–7649.

Deyao Zhu, Jun Chen, Xiaoqian Shen, Xiang Li, and Mohamed Elhoseiny. 2023. MiniGPT-4: Enhancing Vision-Language Understanding with Advanced Large Language Models. In *The Twelfth International Conference on Learning Representations*.

## A Appendix

You are a robot which follows route instructions step-by-step to reach a destination. At every step, you will receive:

1. Route Instruction: the instruction to follow.
2. Current Step: The step number you are currently on in the overall route.
3. Cumulative Distance Traveled: The total distance (in meters) you have moved from the starting point up to your current position.
4. Observations from previous steps (if available), including:
   - Images captured at previous steps.
   - Actions performed at previous steps.
5. Current image: An image showing the robots present view.
6. Possible actions: The set of available actions for this step.

Actions and their definitions:
- Right: Rotates 30 degrees to the right.
- Left: Rotates 30 degrees to the left.
- Move: Moves you forward in your current direction of view.
- Stop: Choose this action when you think you have reached the goal or the end of the navigation path.

Important Notes
- Choose only one action at a time, using only the predefined actions listed in the 'Possible Actions' field.
- The environment is graph-based, meaning movement occurs between discrete nodes rather than continuous space.
- Automatically Turn Towards Node: When you move forward, the camera is automatically adjusted to center on the next node in the graph-based environment. This is handled separately and does not require prediction. Your responsibility is to decide movement and rotation based on the provided inputs.

Your task is to predict the most appropriate next action at each step based on the given information.

Figure 4: System prompt for low-level action space

You are a robot which follows route instructions step-by-step to reach a destination. At every step, you will receive:

1. Route Instruction: the instruction to follow.
2. Current Step: The step number you are currently on in the overall route.
3. Cumulative Distance Traveled: The total distance (in meters) you have moved from the starting point up to your current position.
4. Panorama Images from Previous Steps: If available, these images provide context about where you have been. Use them to understand your past movements and to identify which parts of the current route instruction are most relevant to your current step.
5. Current Panorama Image: A 360-degree panoramic image of your current surroundings. The center of the image represents your current forward direction.
6. Candidate Directions: A list of possible directions to move.
   Each candidate includes:
      - Relative angle: The direction relative to your forward orientation (e.g., '30° left' or '45° right').
      - Distance: The distance (in meters) to the next possible location in that direction.
      - A view (image): What you would see if you move in that direction.
7. STOP Candidate: This is always available and must only be selected when you are certain you have reached the final destination as described in the route instruction.

Your task:
Using the provided inputs, analyze and select the one candidate direction that best matches the route instruction and ensures you stay aligned with the intended path.

Figure 5: System prompt for panoramic action space

With adjustment

Turn Left

Automatically Turn
Towards Node

Move Forward

Without Adjustment

Turn Left

Move Forward

Figure 6: Figure illustrating the `Automatically Turn Towards Node` step.

Figure 7: Figure illustrating the difference between the traditional panoramic approach and our implementation

463

# HiSlang-4.9k: A Benchmark Dataset for Hindi Slang Detection and Identification

**Tanmay Tiwari[1]\*   Vibhu Gupta[1]\*   Manikandan Ravikiran[1]   Rohit Saluja[1,2]**

[1]Indian Institute of Technology Mandi

[2]BharatGen Consortium

{s23107, b22248, erpd2301}@students.iitmandi.ac.in, rohit@iitmandi.ac.in

## Abstract

Slang is an informal register of language, and understanding it is crucial for daily communication. While research on slang detection and identification exists in English (a resource-rich language with abundant data on web), the field remains underexplored in low-resource Indian languages (e.g., Hindi, which has < 1% data on web) due to the lack of comprehensive datasets. Hindi, despite being spoken by over 600 million people worldwide, remains critically underrepresented in Natural Language Processing (NLP) research. In this paper, we introduce HiSlang-4.9k, a dataset containing 4,906 unique sentences, 50% with slang and 50% without slang. HiSlang-4.9k is collected from various resources and is manually annotated with the help of two linguistic experts and eight annotators. We benchmark the performance of state-of-the-art models like BERT, mBERT, IndicBERT, and XLM-RoBERTa on HiSlang-4.9k. We establish benchmarks for slang detection and identification tasks, giving relevant insights into model performance. The IndicBERT model performs the task of slang detection and identification with an F1 score of 0.95 and 0.93, respectively. Additional studies on removing slang and non-slang phrases from sentences during inference highlight models' effectiveness in using the important parts of input for the relevant tasks.

## 1 Introduction

Often used in daily conversation, slang is an informal word or phrase that elicits strong reactions (Coleman, 2012). Lexical flexibility is one of the unique qualities of slang; it lets speakers express ideas creatively across diverse contexts. Slang is an evolving innovation of humans, hence, frequent words that compose slang take on new meanings, and often, whole new terms show up. The evolution makes it difficult for computational

---

\*Equal contribution.



Fig. 1: A slang interpretation example in Hindi where the listener detects (top) and identifies (bottom) the slang. English translation, just for better understanding: He turns every talk upside-down, I don't understand it! (Slang phrase: उलट−पुलट, literal sense: upside-down, slang sense: into confusion).

systems to catch both semantic and pragmatic nuances in slang (Pei et al., 2019). For example, as illustrated in Fig. 1, a listener must infer that the Hindi phrase "उलट−पुलट" (lit. "upside-down") is being used in everyday speech in the Hindi language to convey a sense of disorder, rather than its literal sense (upside-down). This makes it essential for NLP models to first detect Slang (figuring out whether a sentence contains slang) and then identify slang (locating the slang terms), much like how humans interpret non-literal language.

In the last few years, interesting works have been done on processing slang computation-

ally (Dhuliawala et al., 2016; Pei et al., 2019; Sun et al., 2021), but the majority of their work is in English. Lexical resources like SlangNet (Dhuliawala et al., 2016), a network of slang terms, demonstrate how structured slang lexicons can aid downstream NLP tasks in English. There is also a growing interest in using large language models (LLMs) to understand slang better (Sun et al., 2024). The above-mentioned works underscore steady progress in English slang processing; however, similar advances are largely absent for low-resource Indian languages like Hindi, which consists of <1 % of data on web (Q-Success, 2024).

Slang exhibits pronounced semantic divergence between literal and intended meanings. Such context-sensitive and evolving usage makes automatic slang recognition exceedingly difficult without specialized data (Cai et al., 2025). Despite being spoken by more than 600 million people globally, Hindi is still notably underrepresented in NLP research (Thirumala and Ferracane, 2022). To the best of our knowledge, the field of computational slang processing in Hindi is unexplored. Slang in Hindi poses unique challenges that have not been seen in previous English-focused works. Hindi speakers often intermingle English or regional dialect words as slang terms, e.g., आज का क्या सीन है, translation: what is the scene today. The slang term scene/सीन represents plan here. Moreover, Hindi slang sentences often have multiple continuous words or phrases as slang terms, as shown in Figs. 1, 2, 3(b), and Secs. 4.2.1, 4.2.2. Yet, until now, no public dataset or benchmark exists for Hindi slang detection and interpretation. The research gap hinders the development of robust NLP tools for informal Hindi, which are increasingly needed as social media and online content in Indian languages grow. In this work, we address the above-mentioned gaps by introducing HiSlang-4.9k towards benchmarking Hindi slang detection and identification. Our contributions can be summarized as follows:

- We create HiSlang-4.9k: the first dataset for slang detection and identification in Hindi, to the best of our knowledge. HiSlang-4.9k contains 2,453 sentences with slang and 2,453 sentences without slang. Each sentence is manually annotated for slang usage, providing the first resource to study slang detection and identification in Hindi.

- We benchmark several state-of-the-art lan-

guage models for slang detection and identification in the Hindi language. The results provide a comprehensive baseline, with the best fine-tuned model (IndicBERT) achieving F1 scores of 0.95 in detection and 0.93 in identification. We also present additional studies by removing slang and non-slang parts from sentences to identify various challenges in slang detection and identification.

## 2 Related Work

Efforts to build slang-specific lexical resources include SlangNet, which organizes slang terms in a WordNet-like network to better separate senses (Dhuliawala et al., 2016), and SlangSD, a large sentiment dictionary of slang expressions (Wu et al., 2018). These resources highlight the importance of structured lexica in handling informal expressions. Researchers have also explored using contextual embeddings to better model slang. Pretrained models such as BERT (Devlin et al., 2019) and GPT-4 (Achiam et al., 2023) often assign low probabilities to slang terms, making them difficult to detect (Sun et al., 2024). SlangTrack (Anonymous, 2024) addressed this by fine-tuning BERT-large-uncased on English slang data, achieving 87% accuracy in slang detection. Meanwhile LLMs and slangs are recently explored (Sun et al., 2024), using datasets from movie subtitles to test tasks like regional or time-specific slang detection. They showed that GPT-4 performs well in zero-shot settings, while smaller models like BERT can match this performance after fine-tuning on slang-specific data. Besides, slang detection also aligns with broader work in informal language processing. Notable works in this area includes usage of bidirectional LSTMs with POS tagging and character-level convolutional embeddings for sequence labeling of slang, highlighting the syntactic fluidity of slang terms (Pei et al., 2019).

## 3 HiSlang-4.9k Dataset

In this section, we describe the collection and the annotation process of the HiSlang-4.9k dataset, a novel resource designed for slang detection and identification in Hindi. To the best of our knowledge, this is the first dataset on Hindi slang related research. The annotation of HiSlang-4.9k involves eight native annotators (with two annotators labeling each sentence) and two linguistic experts.

तुम्हारी आने की सुनी ही नहीं, दिल गार्डन-गार्डन हो गया मेरा!

●NON-SLANG     ●SLANG     ●NON-SLANG

Fig. 2: Phrase level annotation for slang identification in Hindi highlighting the slang "दिल गार्डन– गार्डन" (lit. "heart garden-garden," Slang sense "felt overjoyed").

## 3.1 Data Source

While the work on slang detection and identification exists for English (Pei et al., 2019), no work exists for Hindi. We first curate 10,000 sentences from diverse sources, including movies scripts and subtitles, linguistic corpora, and online platforms such as social media and discussion forums, following the methodology of (Sun et al., 2024). The subsequent subsections describe the selection and annotation procedure for creating HiSlang-4.9k.

## 3.2 Sentence-level Annotations

The first phase of annotation focuses on classifying sentences based on whether they contain slang or not. Each of the 10,000 sentences is independently reviewed by four annotators who are native Hindi speakers. Each of the 10,000 sentences is labeled by two of the annotators. Annotators are instructed to assign a label based on following:

- **Slang Sentence:** The sentence contains words or phrases that are conventionally used in a non-standard, informal manner in daily conversation.

- **Non-Slang Sentence:** The sentence is fully formal or contains no instances of words or phrases that are conventionally used in a non-standard, informal manner in daily conversation.

With this, we are left with 3,018 sentences identified as slang sentences and the remaining 6,982 sentences as non-slang sentences. We evaluate annotation reliability by computing Cohen's Kappa coefficient (Cohen, 1960). For sentence-level annotations, we observe the Kappa score of 0.97, which happens perhaps because of annotators being native Hindi speakers. As we will see in next subsection, only 2,453 slang sentences are retained in the final dataset based on the inconsistencies in phrase-level annotations. Moreover, 2,453 non-slang sentences from 6,982 are retained to keep the dataset balanced towards the two classes.

**Dataset Statistics**

| Statistic | Value |
| --- | --- |
| Total sentences | 4,906 |
| Slang sentences | 2,453 |
| Non-slang sentences | 2,453 |
| Avg. words per sentence | 15.5 |

Table 1: Key statistics of HiSlang-4.9k.

## 3.3 Phrase-Level Annotations

In the second phase, the 3,018 slang sentences are further subjected to phrase-level annotations, as shown in Fig. 2. The goal is to pinpoint the exact span of slang within each sentence. To reduce bias, four annotators, different from those involved in the sentence-level annotations in the previous subsection, are employed. Each sentence is labeled by two annotators. Precisely, each word in a sentence is labeled as:

- **Slang**: part of a slang phrase.

- **Non-Slang**: not part of any slang phrase.

We observe a Kappa score of 0.94 for phrase-level annotations, suggesting strong consistency between annotators. The high agreement may be attributed to the high frequency of non-slang words in sentences, which made non-slang words easier to agree upon. The annotators' cultural fluency and linguistic intuition as native Hindi speakers familiar with a wide range of slang expressions may also contribute to the high agreement. Although 3,018 sentences are initially labeled as containing slang at the sentence-level, only 2,453 are retained by the two experts in the final dataset based on inconsistencies in the annotations. The experts also merge the annotations with slight differences. The final dataset consists of 4,906 sentences, with 50% of sentences with slang and the remaining 50% without slang (sampled from original non-slang sentences to keep the two classes balanced, as discussed in the previous subsection). The key statistics of the HiSlang-4.9k dataset are summarized in Table 1. The next paragraphs present more detailed insights on slang sentences in HiSlang-4.9k.

We now analyze the phrase-level annotations in the HiSlang-4.9k dataset to understand its properties. Some interesting distributional patterns emerge for slang usage; Fig. 3a illustrates the distribution of the position of slang words within

(a) Distribution of slang words' position in the sentences.

(b) Number of slang words per sentence.

(c) POS distribution of slang words

Fig. 3: Analysis of Slang Sentences in HiSlang-4.9k. the slang sentences, with positions normalized from 0 to 1. As shown, the slang words tend to exist toward sentence boundaries, appearing more frequently near the start or the end of a sentence than in the middle. Therefore, slang often serves as an opener or closer in informal Hindi statements, a pattern that detection and identification models can exploit by incorporating positional features.

Fig. 3b depicts the count of slang words per sentence in the subset of samples with slang. The bell-shaped curve displays the variety of slangs in HiSlang-4.9k, with the majority of sentences having three slang words.

Fig. 3c presents the part-of-speech (POS) distribution of slang words. Nouns (NOUN) and verbs (VERB) dominate Hindi slang usage, followed by smaller proportions of adpositions (ADP), auxil-

iary verbs (AUX), and adjectives (ADJ). This POS usage skew reveals a lexical preference for using slang in descriptive and referential roles, reflecting how informal Hindi relies on creative nouns and verbs. Such trends provide valuable cues for slang detection and identification; knowing that slang is often a noun or verb can guide models to focus on these word classes when distinguishing slang from the standard lexicon. Such insights, models trained on Indian languages (see Sec. 2), and the recent work on Indic MCQs (Ravikiran et al., 2025) help us select the pretrained models and fine-tuning strategies we discuss in the next section.

## 4 Experiments and Results

In this section, we present the results of different Indic language models for the detection and identification tasks. The tasks are defined by Sun et al. (2024). Slang detection refers to the classification task where a model determines whether a given sentence contains at least one instance of slang usage. Slang identification is a more fine-grained task in which models perform *phrase-level tagging* to pinpoint the exact words or spans within a sentence that constitute a slang phrase.

All experiments are performed using transformer-based architecture, specifically BERT (Devlin et al., 2019), mBERT (Devlin et al., 2019), IndicBERT (Kakwani et al., 2020) and XLM-RoBERTa (Conneau et al., 2020) as used in a recent work on Indic MCQ difficulty estimation (Ravikiran et al., 2025). Another reason for using the abovementioned models is that all models except the BERT are pretrained on Hindi data. The fine-tuning strategy for slang detection and identification are the same as classification and Named Entity Recognition (NER) tasks based on the dataset analysis discussed in the previous section and protocols defined by (Ravikiran et al., 2025) and (Wolf et al., 2020). We evaluate the performance by means of two configurations: (a) Complete model fine-tuning, whereby the final detection/identification layer is added to the model and the entire model is fine-tuned on pretrained weights; (b) Last layer fine-tuning, whereby the last layer added is fine-tuned with the frozen pretrained weights.

We split the data into train:test ratio of 80:20. Each model is assessed using the performance metrics of Precision (P), Recall (R) and F1-Score (F1).

### (a) Slang Detection

| Model | P | R | F1 |
|---|---|---|---|
| IndicBERT | 0.9379 | **0.9554** | **0.9466** |
| XLM-RoBERTa | **0.9535** | 0.9145 | 0.9336 |
| mBERT | 0.9358 | 0.9220 | 0.9289 |
| BERT | 0.8606 | 0.8022 | 0.8302 |
| XLM-RoBERTa† | 0.7655 | 0.8736 | 0.8167 |
| mBERT† | 0.7255 | 0.6877 | 0.7061 |
| IndicBERT† | 0.6616 | 0.7104 | 0.6851 |
| BERT† | 0.7252 | 0.5985 | 0.6551 |

### (b) Slang Identification

| Model | P | R | F1 |
|---|---|---|---|
| IndicBERT | 0.9221 | **0.9332** | **0.9276** |
| XLM-RoBERTa | **0.9305** | 0.9105 | 0.9204 |
| mBERT | 0.9018 | 0.9093 | 0.9055 |
| BERT | 0.8761 | 0.8938 | 0.8849 |
| IndicBERT† | 0.8018 | 0.4345 | 0.5634 |
| mBERT† | 0.7544 | 0.3555 | 0.4823 |
| XLM-RoBERTa† | 0.4321 | 0.1897 | 0.2631 |
| BERT† | 0.2627 | 0.0680 | 0.1078 |

Table 2: Results of BERT-based models on (a) slang detection and (b) slang identification. Metrics: Precision (P), Recall (R), F1. Models: BERT (Devlin et al., 2019), mBERT (Devlin et al., 2019), IndicBERT (Kakwani et al., 2020), XLM-RoBERTa (Conneau et al., 2020). Models marked with † are *frozen*: only the final layer is trained.

## 4.1 Results

Table 2 presents the performance of various models on the two tasks: (a) slang detection and (b) slang identification. Across both tasks, we observe a consistent trend where models whose pretraining included Hindi data, namely, IndicBERT, XLM-RoBERTa and mBERT outperform the English-only model BERT. This indicates that familiarity with Indian linguistic patterns and vocabulary substantially improves slang processing capabilities. For slang detection results shown in rows 1-4 of Table 2 (a), IndicBERT achieves the highest F1 score of 0.9466, closely followed by XLM-RoBERTa at 0.9336. IndicBERT performs best due to pretraining of the dataset in Indian languages, while XLM-RoBERTa and mBERT include a mix of Indic and non-Indic languages. XLM-RoBERTa has higher precision compared to IndicBERT, possibly due to the involvement of English transliterations in the slang data (see example on usage of slang term scene/सीन in Sec 1). In contrast, BERT lags at 0.8302, suggesting that it fails to model Hindi

slang effectively due to its English-centric pretraining. The lower F1-Scores with last-layer finetuning in rows 5-8 of Table 2 (a) with respect to rows 1-4 suggest that merely updating the last layer is not as good as updating all the layers of the models. Hence, we can conclude that although initializing the weights with Indian data helps, slang detection is a complex task and hence requires the transformation of all the weights in the models.

The results of the slang identification task with complete model fine-tuning are shown in rows 1-4 of Table 2 (b). Similar to the detection results, IndicBERT achieves the highest performance with an F1-score of 0.9276, followed by XLM-RoBERTa at 0.9204 and mBERT at 0.9055. BERT shows lower performance with an F1-score of 0.8849. The higher gaps between rows 1-4 and rows 5-8 (last-layer fine-tuning) of Table 2 (b), compared to the detection task (previous paragraph), show that slang identification, being a finer task than detection, is even harder with the single last-layer fine-tuning.

The results validate the quality and the complexity of the HiSlang-4.9k dataset. Models trained and evaluated on it exhibit meaningful performance differences that align with their expected linguistic capabilities. Moreover, HiSlang-4.9k appears to be a reliable benchmark for both sentence classification and phrase-level tagging in informal language processing for Hindi.

## 4.2 Qualitative Analysis

In this section, we conduct a qualitative analysis of the results obtained by IndicBERT (Kakwani et al., 2020) fine-tuned on the HiSlang-4.9k.

### 4.2.1 Slang Detection

For the qualitative analysis of the slang detection task, we observe various patterns in the model's decision-making. The labels and predictions are marked as ✓ for informal slang sentences, while for a formal sentence without slang, the same are marked as ✗. Below, we present representative examples from both success and failure cases, along with a brief analysis of each.

**Success Cases:**

- Sentence: पूरी तैयारी के बावजूद फाइनल मैच में हारने से टीम की शान का बर्था बन गया। (Despite all the preparation, losing the final match turned the team's pride into a mash.)

468

Ground Truth: ✓     Prediction: ✓

The model correctly identifies this sentence as slang due to the use of the informal phrase बर्था बन गया (literal sense: "turned into a mash", slang sense "utterly destroyed"). The context and the informal phrase, in addition to the position of the slang phrase and usage of verbs and nouns in the slang phrase (see Figs. 3a, 3c and Sec. 3.3), likely helped the model capture the slang intent.

- Sentence: उन्होंने कार्यालय में सभी लेख्य सही समय पर जमा किए। (He submitted all the documents correctly on time at the office.)
Ground Truth: ✗     Prediction: ✗
The model correctly labels this as a non-slang sentence because it uses formal vocabulary and a clear declarative structure, with no informal expressions to suggest slang.

**Failure Cases:**

- Sentence: उसने छोटा–मोटा काम अपनी बहन से करवाया और खुद फोन पर खेलता रहा। (He had his sister do small-fat work, while he himself kept playing on his phone.)
Ground Truth: ✓     Prediction: ✗
The model labels this as a non-slang sentence because "छोटा–मोटा काम" (literal sense: "small–fat work", slang sense: "a little bit of work") is composed of some of the terms ("छोटा काम"/small work) which jointly have a similar meaning to the slang sense.

- Sentence: उसने अंगूठा–छाप आदमी को प्रोजेक्ट सौंप दिया। (He assigned the project to the thumb print person.)
Ground Truth: ✓     Prediction: ✗
The slang word "अंगूठा–छाप" (literal sense: "thumb print", slang sense: "illiterate") contributes to an informal tone, but the model fails to identify it as containing slang. This may be due to the formal tone of the rest of the sentence overshadowing the slang word, leading to misclassification.

These examples highlight that while IndicBERT performs well in cases with explicit or contextual slang cues, it sometimes struggles with slang terms having similar meaning to slang sense, and may also under-detect single-word slang in more formal constructions.

### 4.2.2 Slang Identification

For the qualitative analysis of the slang identification task, which involves predicting slang words from the sentence, we observe a range of predictions across different kinds of slang sentences. One pattern of error is that the model misses words that are at the boundary of the slang phrase. Words such as मैं, से, है, दिया are excluded from the predicted span, even though they are part of the ground truth and contribute significantly to the interpretability of the slang. Representative examples of such errors include:

- Sentence: उसने अपने परिवार की इज्जत को मिट्टी में मिला दिया। (He completely mixed his family's honor into the soil.)
Ground Truth: मिट्टी में मिला दिया (literal sense: "mixing into the soil", slang sense: "utterly destroy or humiliate")
Prediction: मिट्टी में मिला

- Sentence: पहले इंटरव्यू में पास होते ही मुझे चांदी हो जाना महसूस हुआ। (As soon as I passed the first interview, I felt I became silver.)
Ground Truth: चांदी हो जाना (literal sense: "become silver," slang sense: "got lucky")
Prediction: चांदी हो

The above-mentioned failure cases are possibly due to boundary terms being less frequent parts-of-speech (POS) terms as shown in Fig. 3c.

Contrary to the above example, in the case of a single-word slang, the model is generally able to identify the slang term, but sometimes includes extra words from the surrounding context. This behavior often leads to over-extended predicted spans. Such errors are also possibly due to over-fitting on highly frequent multi-word slangs (see Fig. 3b). Illustrative examples are:

- Sentence: टीम मीटिंग में ढक्कन ने ऐसा सुझाव दिया कि सबका ध्यान उसकी बेवकूफी पर चला गया। (In the team meeting, the lid made such a suggestion that everyone's attention shifted to his foolishness.)
Ground Truth: ढक्कन (literal sense: "lid", slang sense: "fool")
Prediction: ढक्कन ने

- Sentence: दोगला इंसान हमेशा अपनी बातों और कामों में विरोधाभास रखता है। (A person with two necks always keeps contradictions between their words and their actions.)

| Model | Exp. 1 | Exp. 2 | Exp. 3 |
|---|---|---|---|
| BERT | 0.8664 | 0.8914 | 0.8103 |
| mBERT | 0.8861 | 0.9155 | 0.9117 |
| IndicBERT | **0.9237** | 0.9287 | 0.9308 |
| XLM-RoBERTa | **0.9237** | 0.9330 | 0.9332 |

Table 3: F1 scores (slang detection) across three additional experiments: Exp. 1—non-slang sentences only; Exp. 2—slang removed from slang sentences; Exp. 3—isolated slang phrases.

Ground Truth: दोगला (literal sense: "two necks", slang sense "hypocrite")
Prediction: दोगला इंसान

Finally, there are several correct predictions where the slang is identified with precise boundaries, indicating a successful understanding by the model:

- **Sentence:** पोपट बना दिया उसने अपनी झूठी कहानियों से मुझे। (He made a parrot of me with his false stories.)
  **Ground Truth:** पोपट बना दिया (literal sense: "made a parrot", slang sense: "made a fool")
  **Prediction:** पोपट बना दिया

- **Sentence:** ज़्यादा उछल रहा है तु आजकल। (You've been jumping too much these days.)
  **Ground Truth:** ज़्यादा उछल रहा है (literal sense: "jumping too much", slang sense: "reckless")
  **Prediction:** ज़्यादा उछल रहा है

### 4.3 Additional Studies

Table 3 summarizes the F1-score of each model across three additional experiments designed to probe their ability to distinguish slang from non-slang under increasingly challenging conditions. The three experiments are performed for the slang detection task.

In Experiment 1, where fine-tuned models are applied to the 2,453 non-slang sentences, IndicBERT and XLM-RoBERTa achieve the highest F1-score (92.37%), whereas mBERT and BERT lag at 88.61% and 86.64% respectively, frequently mislabeling non-slang words as slang. XLM-RoBERTa, IndicBERT (92.37%) and mBERT (88.61%) both outperformed BERT, indicating that Hindi-aware models work even when explicit slang cues are absent.

In Experiment 2, slang phrases were removed from the 2,453 slang sentences that originally con-

tained them, so models had to rely solely on context. The labels are also modified from slang to non-slang in this case. XLM-RoBERTa again led the field at 93.30% F1-score. IndicBERT (92.87%) and mBERT (91.55%) both show good performance, confirming that their Hindi-aware pretraining helps. BERT (89.14%) remains the least reliable.

In Experiment 3, models are evaluated on isolated slang phrases with no surrounding context. XLM-RoBERTa again performs the best (93.32%), correctly identifying most slang expressions in isolation. IndicBERT (93.08%) and mBERT (91.17%) follow the performance closely. In contrast, BERT's performance drops to 81.03%, underscoring its difficulty in recognizing slang without additional contextual information.

These results validate that fine-tuned XLM-RoBERTa and IndicBERT are effective for Hindi slang detection under varying and extreme conditions.

## 5 Conclusion

In this work, we introduce a high-quality dataset, HiSlang-4.9k, for slang detection and identification in the Hindi language. Recognizing the growth of informal online communication, especially involving slang, the dataset addresses a gap in existing language resources for Indian languages. The corpus comprises 4,906 manually annotated sentences, sourced from the real-world text. We employed carefully designed annotation guidelines and a rigorous validation process to ensure a high-quality dataset. The dataset includes both slang and non-slang sentences, with diverse sentence structures that feature slang words in varying positions and contexts. To assess the utility of HiSlang-4.9k, we benchmark multiple transformer-based model architectures. Our experiments demonstrate that Hindi-language pretrained models (e.g., IndicBERT, XLM-RoBERTa, mBERT) fine-tuned on our data significantly outperform the English-only model (BERT), highlighting the importance of language-aware training.

Overall, we believe that HiSlang-4.9k, along with the benchmarks established in this work, can serve as a valuable foundation for future research in informal-language processing for the Hindi language. We will release the dataset and baseline implementations to encourage further exploration in

this direction.

## Limitations

The study makes progress in handling informal Hindi, but it also has limitations. Annotating slang depends on people's views; a phrase interpreted by one person as slang might be a non-slang term for the other. This subjectivity arises because speakers come from varied backgrounds and use words in different settings. We try to mitigate these gaps by giving clear instructions and using multiple annotators for each example. Even so, some variation in how slang is interpreted may still appear across our data. Also, the dataset might not fully represent the variety of Hindi slang. Most of our data comes from social media and online forums, reflecting mainly the language used by younger people familiar with the internet. Because of this, slang from other communities, dialects, or regions may not be well-covered. This limitation means our models might struggle with slang terms common in these underrepresented groups.

## Ethics Statement

We acknowledge that slang is inherently subjective and can be sensitive in certain contexts, especially in informal speech where meanings may vary widely across different communities and generations. All annotations and analyses were performed by native Hindi speakers, and the data was sourced from publicly available content such as social media posts and discussion forums. We ensured that all contributors to data annotation and analysis participated voluntarily and were informed of the research goals. No personally identifiable or sensitive information was collected or shared. We understand that certain slang terms might carry connotations that could be considered offensive or inappropriate in some contexts. We therefore encourage users of the HiSlang-4.9k dataset to apply cultural sensitivity and appropriate disclaimers when deploying models or sharing results derived from this dataset. Our aim is solely to advance the understanding of informal Hindi language in NLP research and to promote inclusive and responsible use of linguistic resources.

## Acknowledgements

## References

Josh Achiam, Steven Adler, Sandhini Agarwal, Lama Ahmad, Ilge Akkaya, Florencia Leoni Aleman, Diogo Almeida, Janko Altenschmidt, Sam Altman, Shyamal Anadkat, et al. 2023. Gpt-4 technical report. *arXiv preprint arXiv:2303.08774*.

Anonymous. 2024. Slang or not? exploring nlp techniques for slang detection using the slangtrack dataset. arXiv preprint arXiv:2401.00001. Submitted to ACL ARR 2024.

Jinyu Cai, Yusei Ishimizu, Mingyue Zhang, Munan Li, Jialong Li, and Kenji Tei. 2025. Simulation of language evolution under regulated social media platforms: A synergistic approach of large language models and genetic algorithms. *arXiv preprint arXiv:2502.19193*.

Jacob Cohen. 1960. A coefficient of agreement for nominal scales. *Educational and Psychological Measurement*, 20(1):37–46.

Julie Coleman. 2012. *The life of slang*. Oxford University Press, USA.

Alexis Conneau, Kartikay Khandelwal, Naman Goyal, Vishrav Chaudhary, Guillaume Wenzek, Francisco Guzmán, Edouard Grave, Myle Ott, Luke Zettlemoyer, and Veselin Stoyanov. 2020. Unsupervised cross-lingual representation learning at scale. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 8440–8451, Online. Association for Computational Linguistics.

Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. 2019. BERT: Pre-training of deep bidirectional transformers for language understanding. In *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers)*, pages 4171–4186, Minneapolis, Minnesota. Association for Computational Linguistics.

Shehzaad Dhuliawala, Diptesh Kanojia, and Pushpak Bhattacharyya. 2016. SlangNet: A WordNet like resource for English slang. In *Proceedings of the Tenth International Conference on Language Resources and Evaluation (LREC'16)*, pages 4329–4332, Portorož, Slovenia. European Language Resources Association (ELRA).

[1]BharatGen

Divyanshu Kakwani, Anoop Kunchukuttan, Satish Golla, Gokul N.C., Avik Bhattacharyya, Mitesh M. Khapra, and Pratyush Kumar. 2020. IndicNLPSuite: Monolingual corpora, evaluation benchmarks and pre-trained multilingual language models for Indian languages. In *Findings of the Association for Computational Linguistics: EMNLP 2020*, pages 4948–4961, Online. Association for Computational Linguistics.

Zhengqi Pei, Zhewei Sun, and Yang Xu. 2019. Slang detection and identification. In *Proceedings of the 23rd Conference on Computational Natural Language Learning (CoNLL)*, pages 881–889, Hong Kong, China. Association for Computational Linguistics.

Q-Success. 2024. Usage of hindi broken down by content management systems. Accessed: 2025-06-08.

Manikandan Ravikiran, Siddharth Vohra, Rajat Verma, Rohit Saluja, and Arnav Bhavsar. 2025. Teemil: Towards educational mcq difficulty estimation in indic languages. In *Proceedings of the 31st International Conference on Computational Linguistics*, pages 2085–2099.

Zhewei Sun, Qian Hu, Rahul Gupta, Richard Zemel, and Yang Xu. 2024. Toward Informal Language Processing: Knowledge of Slang in Large Language Models. In *Proceedings of NAACL□HLT 2024 (Long Papers)*, pages 1683–1701, Mexico City, Mexico. Association for Computational Linguistics.

Zhewei Sun, Richard Zemel, and Yang Xu. 2021. A computational framework for slang generation. *Transactions of the Association for Computational Linguistics*, 9:462–478.

Adhitya Thirumala and Elisa Ferracane. 2022. Extractive question answering on queries in hindi and tamil.

Thomas Wolf, Lysandre Debut, Victor Sanh, Julien Chaumond, Clement Delangue, Anthony Moi, Pierric Cistac, Tim Rault, Remi Louf, Morgan Funtowicz, Joe Davison, Sam Shleifer, Patrick von Platen, Clara Ma, Yacine Jernite, Julien Plu, Canwen Xu, Teven Le Scao, Sylvain Gugger, Mariama Drame, Quentin Lhoest, and Alexander Rush. 2020. Transformers: State-of-the-art natural language processing. In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing: System Demonstrations*, pages 38–45, Online. Association for Computational Linguistics.

Liang Wu, Fred Morstatter, and Huan Liu. 2018. Slangsd: building, expanding and using a sentiment dictionary of slang words for short-text sentiment classification. *Language Resources and Evaluation*, 52:839–852.

# Enhancing Arabic Retrieval Augmented Generation through Language Processing

**Shadi Saleh[1], Belkacem EL Jattari[2], Layth Oud[1], Maryam Alblooshi[1],**
**Bouthaina Lakhdari[2], Ali Alnaqbi[1],**
[1]Prime Technologies, Czechia,
[2]Mohammed bin Zayed University for Humanities, Humanities and Social Sciences
**Correspondence:** shadi@prime-itech.com

## Abstract

We investigate the impact of text preprocessing on Arabic Information Retrieval (IR) systems and, consequently, on the quality of Retrieval-Augmented Generative (RAG) systems. Our work focuses on academic content in Arabic. We analyze how the IR performance affects the quality of RAG systems in answering users' questions on various academic topics. Our findings indicate that the performance of an IR system is significantly influenced by the quality of Optical Character Recognition (OCR) applied to PDF files. We employ a state-of-the-art deep learning-based OCR solution to create our IR index. Eventually, this IR index is used to generate a context-window for the generative model that is employed in chat assistant to answer questions in the scientific domain. We introduce a benchmark dataset for the IR system, comprising 170 Arabic queries and IR relevance assessment with numerous query-document judgment pairs. Our results demonstrate that advanced text preprocessing can lead to an increase of 8 points in terms of $P@5$ of the IR mode, an increase of 11% in the accuracy of the answering system, and up to 95% of correct citations compared to our baseline system.

## 1 Introduction

The Arabic language presents unique challenges for Natural Language Processing (NLP) applications, particularly in Information Retrieval (IR), Question Answering (QA), and Optical Character Recognition (OCR). These challenges stem from the script's complexity, including its cursive nature, contextual letter forms, and bidirectional text orientation. Additionally, the scarcity of high-quality annotated datasets for Arabic further exacerbates these issues. In recent years, Retrieval-Augmented Generation (RAG) systems have emerged as a promising approach to enhance generative models like GPT by grounding their outputs in relevant retrieved documents. This paradigm has proven effective in reducing hallucinations and improving factual accuracy by providing contextually relevant information to generative models.

This paper investigates the impact of advanced Arabic preprocessing techniques on IR systems and their downstream influence on RAG performance. Specifically, we focus on academic content from Arabic scientific papers published in PDF format where accurate OCR and IR are critical for generating high-quality responses. By employing a state-of-the-art deep learning-based OCR solution and optimizing the IR pipeline, we demonstrate significant improvements in both retrieval precision and the relevance of generated answers. Our contributions include a comprehensive evaluation of OCR tools for Arabic PDF processing, and their impact on IR performance, and the development of an Arabic RAG system capable of answering academic queries with enhanced relevance and citation accuracy.

## 2 Related Work

Recently, RAG systems have attracted significant attention from researchers tackling various downstream tasks, particularly in the development of customized chatbots. RAG combines the strengths of retrieval and generative models, enabling chatbots to access grounding knowledge-base and produce contextually relevant responses. This approach has proven effective in creating chatbots that are tailored to in-domain data, by building a retrieval system that searches within in-domain data index and provides context to generative model which has shown to significantly reduces hallucination and provides timely updated information. The ability of generative models to synthesize information from diverse sources showed recently that such models can be used to achieve comparable results to state-of-the-art methods in various NLP tasks (Brown et al., 2020), such as translation, translation-quality evaluation (Kocmi and Federmann, 2023), and

question-answering, usually applying generative models in down-stream tasks requires a few-shot learning approach that can be applied post model training via system prompt.

Sadek et al. (2012) presented one of the first Arabic question-answering system designed to answer *why* and *how* questions. Their approach relies on the discourse structure of Arabic texts to automatically find answers. This method uses Rhetorical Structure Theory (RST), which has been proven effective in NLP applications.

Studies indicated that off-the-shelf OCR tools, which perform considerably well on English, struggle with Arabic due to the script's complexity and the scarcity of high-quality labeled datasets for training purposes. Such models often utilize Hidden Markov Model for sequence modeling (Bunke et al., 1995) . In order to achieve reasonable performance, more advanced method including leveraging language modeling and deep learning approaches (Bhatia et al., 2024).

IR-RAG @ SIGIR24 is a dedicated workshop in SIGIR that emphasizes on the critical rule of IR systems as an internal component of RAG system (Petroni et al., 2024). Multiple submissions argued that the effectiveness of RAG systems heavily relies on the quality of retrieved documents, as poor-quality or irrelevant sources can lead to misleading outputs, and called for a further exploration of robust retrieval mechanisms to enhance RAG capabilities [1].

## 3 Experiment Settings

### 3.1 Data

#### 3.1.1 Queries

We leverage multiple data sources to enhance our analysis. Our dataset comprises 12,000 PDF files of Arabic journal papers sourced from the Shamra Academia portal[2]. This portal is accessible online and indexed by major search engines such as Google[3] and Bing[4]. These search engines provide a search console for website owners, offering insights into visitor engagement, clicks, and search queries.

We analyzed the query logs from the past six months, focusing on the top 100 most frequently searched queries each month. This approach

yielded a total of 600 queries covering various academic research topics. To ensure a diverse and comprehensive set of queries, three experts (native Arabic speakers and academic researchers) selected 170 queries from this pool. We split these queries into 100 queries for training and 70 queries for testing. The split takes into account preserving the distribution of query categories. Table 1 shows the distribution of query categories on both the training and test set. Queries were manually segmented into 8 categories, based on the field of the study related to researches that are relevant to each query.

| Category | train queries | test queries |
|---|---|---|
| Economy | 14 | 10 |
| History | 11 | 8 |
| Engineering | 15 | 11 |
| Agriculture | 14 | 10 |
| Science | 6 | 4 |
| Math | 4 | 3 |
| Medicine | 10 | 7 |
| Literature | 26 | 18 |

Table 1: Distribution of query categories on both training and test sets

Table 2 shows a sample of query text and their categories. The query text is what users type in the search engine to find information that meets their information need. It is important to mention that query text is not necessarily a good representation of searcher's information need. When searchers need to find information about a specific topic, they predict how an ideal document will look like, and what might be the main keywords in that document (e.g. the document title, headers and sub-headers), and they expect the IR system to find those documents for them (El Zein and da Costa Pereira, 2022).

| Category | Query |
|---|---|
| Literature | إدارة التغيير التربوي |
| Engineering | إدارة البيانات السحابية |
| Medicine | التصلب اللويحي |
| History | منهج التأليف في السنة النبوية |

Table 2: Samples of queries chosen and proofread by expert annotators

---

[1] https://ceur-ws.org/Vol-3784/

[2] https://shamra-academia.com

[3] https://google.com

[4] https://bing.com

### 3.1.2 Documents

We have access to 12,000 documents in the portal. These documents are academic papers published in various open-access journals with which the portal has agreements.

Each document in the portal has the following entries (in addition to the PDF file): title, abstract, keywords and the main content. Those fields were added by authors and verified by the portal editors. However, there is still large amount of text in the PDF files needs to be extracted and added to the index so it can be searchable (the main content). To extract the text from PDF files, we experiment with two methods: The first method is based on an open-source python library PyPDF [5], and the second method is based on a more advanced method that utilizes Deep Learning (Surya)[6].

Table 3 presents a sample sentence extracted from a PDF file after conversion to text using two different libraries. The PyPDF library consistently exhibits issues such as merging two words by eliminating the space between them or omitting letters from words. This results in significant information loss. In the sample shown, out of 15 words, PyPDF correctly extracted only 6 words. In contrast, the Surya library accurately extracted the entire sentence without any error. These findings were consistently observed across other documents as well. The primary reasons for PyPDF's poor performance with Arabic text are the complexity of the Arabic script and the bidirectional nature of the language. Arabic script includes cursive writing, contextual letter forms, and diacritics, which are challenging for OCR systems primarily designed for Latin scripts like English. Additionally, Arabic is written from right to left, adding another layer of complexity that PyPDF may not handle effectively.

We observed some issues with Surya OCR system as well when handling Arabic terms that are in-domain, e.g. scientific terms in the medical and engineering domains. To further improve the output of the Surya OCR system and restore information loss, we leverage an LLM as an auto-correction system, where we ask the LLM to improve the output of OCR using the following prompt:

> *You are an expert copyeditor specializing in academic and scientific Arabic texts. Your task is to correct errors in a given OCR-scanned paragraph, consider*

*the following:*
***Text Source:*** *The input is from an Optical Character Recognition (OCR) system. Therefore, you must identify and correct common OCR errors, which include spelling mistakes, garbled words, and incorrect character recognition.*
***Content Type:*** *The text is academic and contains specific scientific terminology. You must preserve all original technical terms and academic phrases without alteration or simplification.*
***Correction Scope:*** *Your corrections should focus exclusively on spelling, grammar, and reconstructing garbled words to make the text fluent and accurate. Do not rewrite, rephrase, or change the intended meaning of the original content.*
***Output Format:*** *Provide ONLY the fully corrected Arabic text. Do not include any introductory phrases, explanations, or commentary.*
***Text to be corrected:***
*{{paragraph_input}}*

From observations with preliminary prompt designs, it was noted that the LLM occasionally tended to rephrase the input text or introduce additional words, altering the original meaning. To remedy this and address common OCR-related errors such as character misrecognition and garbled words, the proposed LLM prompt explicitly identifies the text as OCR output and is designed to instruct the model to preserve domain-specific terminology while strictly avoiding any rephrasing or insertion of extra words.

### 3.1.3 Annotations

To evaluate the performance of our IR system, we conducted a human evaluation using the open-source tool: *relevation*[7]. Relevation is a web application in which human judges can evaluate the relevance of the retrieved documents.

We asked three human judges, all are academic researchers and native Arabic speakers, to perform the evaluation. Each judge was presented with a query and a corresponding document in PDF format. Alongside the query, we provided the information need behind it, and the judges' task was to determine whether the document satisfied the

---

[5] https://github.com/py-pdf/pypdf
[6] https://github.com/VikParuchuri/surya

[7] https://github.com/ielab/relevation

| OCR System | Sentence |
|---|---|
| PyPDF | تقييمحدوث الإستنشاق الري لدمرضى العناية الحرحة الموضوع ليم التغذية المعوة بالانبوب الأفي المعدي |
| Surya | تقييم حدوث الإستنشاق الرئوي لدى مرضى العناية الحرجة الموضوع لهم التغذية المعوّية بالانبوب الأنفى المعدى |

Table 3: Samples of sentence taken from a PDF file of a paper published in a medical journal. The sentence is taken after converting the PDF file into text using two OCR librares, PyPDF

user's need. We employed a ranked evaluation approach, where each document can be annotated as highly relevant, relevant, somewhat relevant, or irrelevant based on how well it addresses the user's information need. At the end of the evaluation, we analyzed the agreement rate among the judges. Each judge was shown 30 document-query pairs that had been previously judged by another judge, and they were asked to re-evaluate these pairs. This process provided insights into the reliability of the judgments. We then binarized the judgments into a 0-1 scale, where irrelevant and somewhat relevant were coded as 0, and relevant and highly relevant as 1. Before binarization, the agreement rate was 70%, which increased to 78% after binarization. This increase is understandable given the varying degrees of relevance perceived by the judges, with the highest disagreement occurring between the somewhat relevant and relevant degrees.

### 3.2 Indexing

To build the IR system that is used to inject context to the LLM prompt, we employed Elasticsearch to retrieve the top 100 most pertinent studies for each query in our dataset. We constructed two distinct indexes for this purpose: one utilizing PyPDF and the other based on Surya. The search filter was designed to optimize the relevance of the results, and include all possible information in each document. The Elasticsearch query schema used is as follows:

```
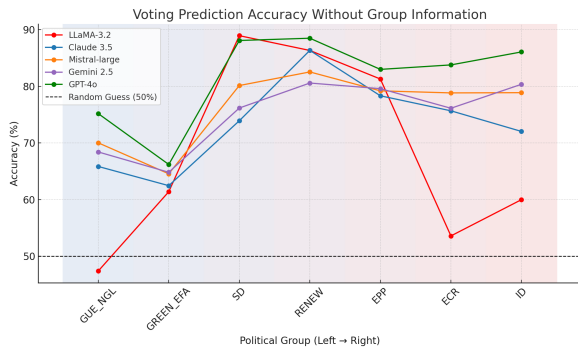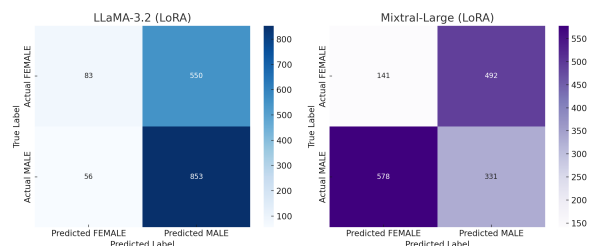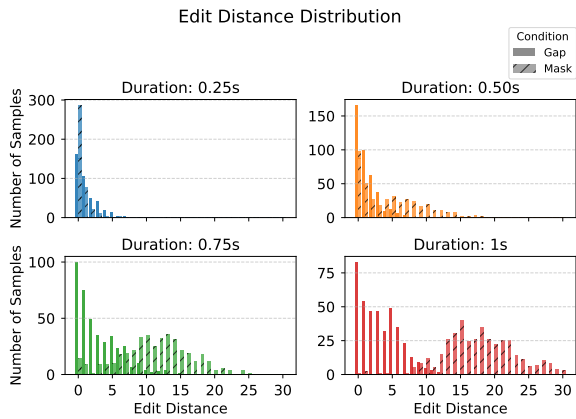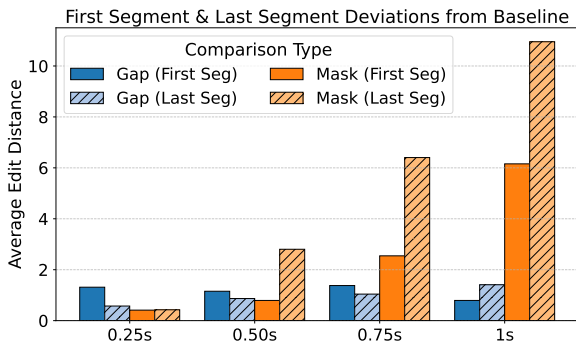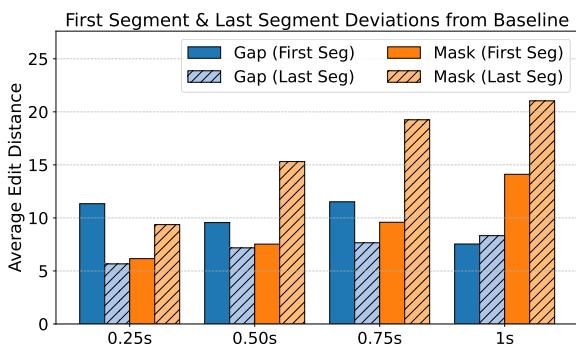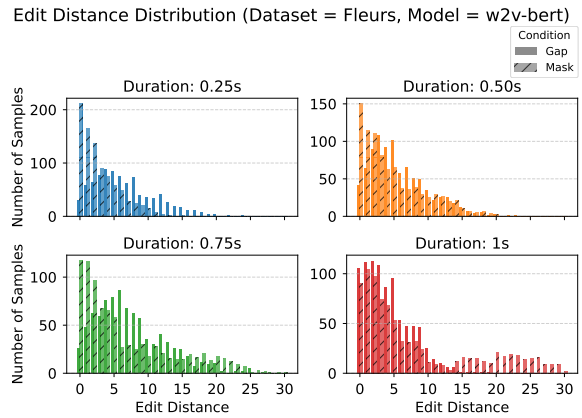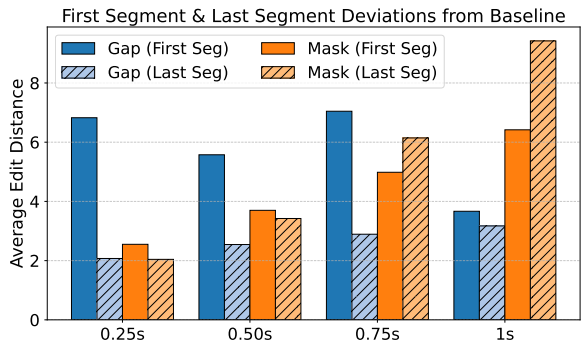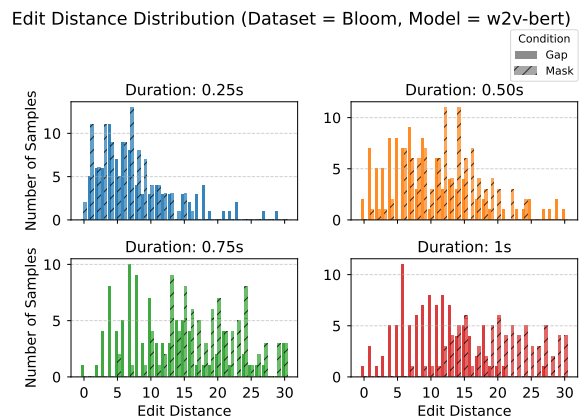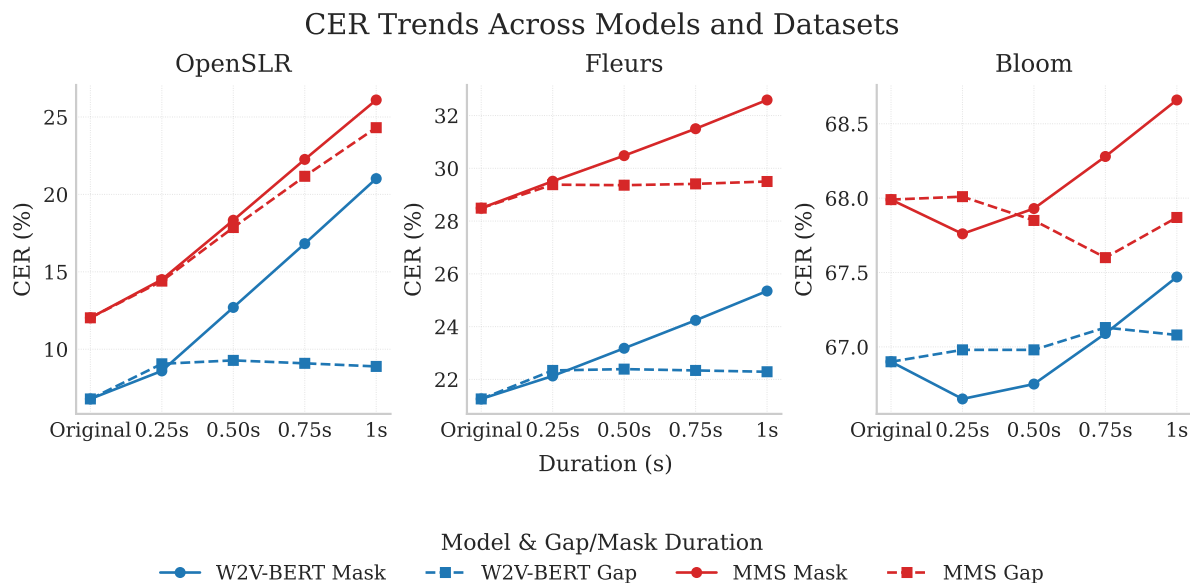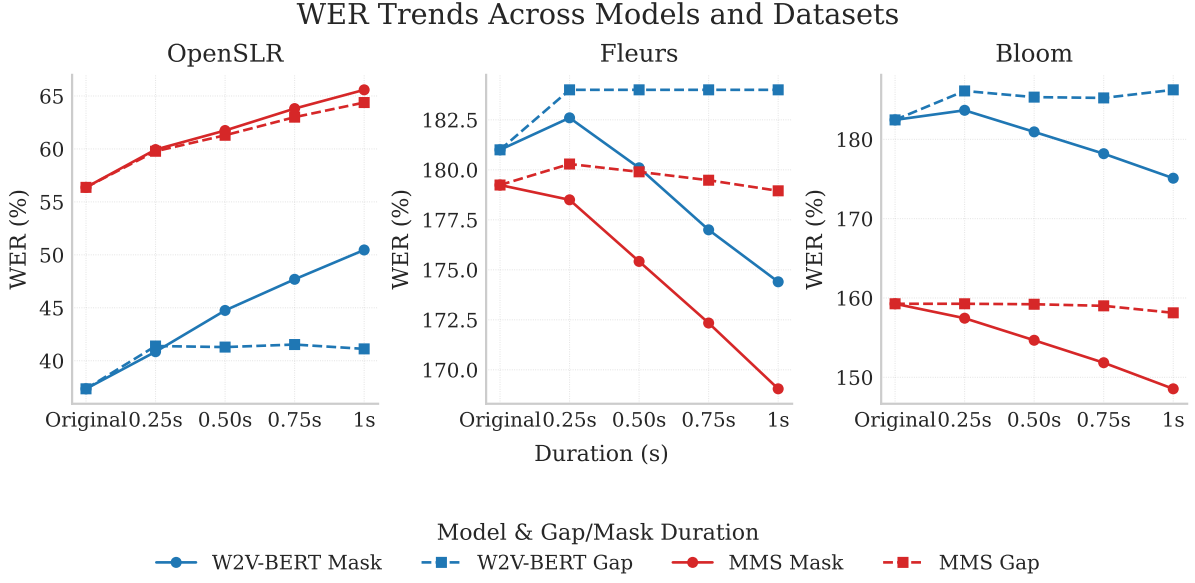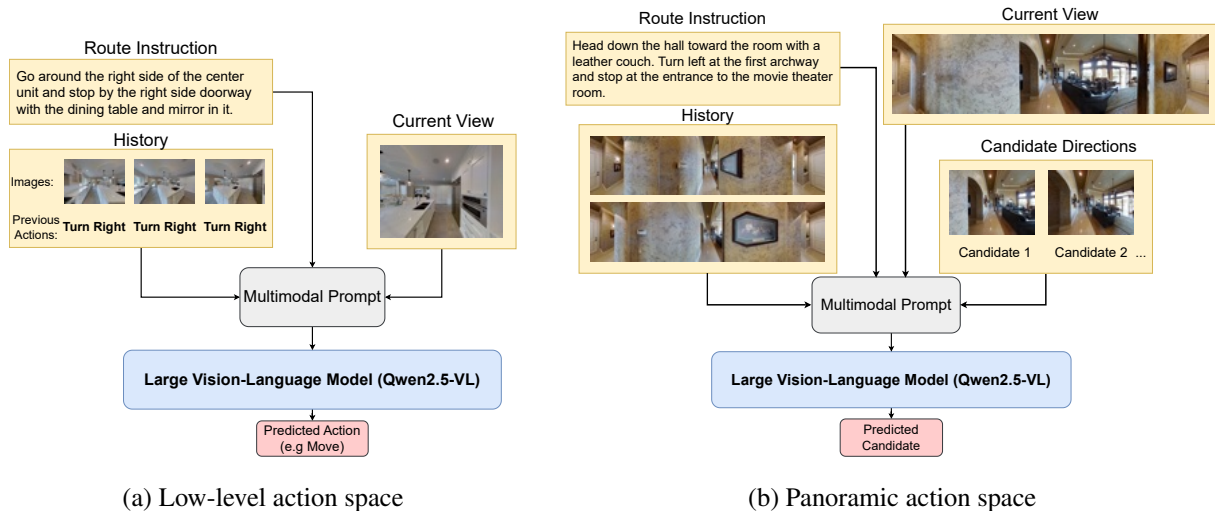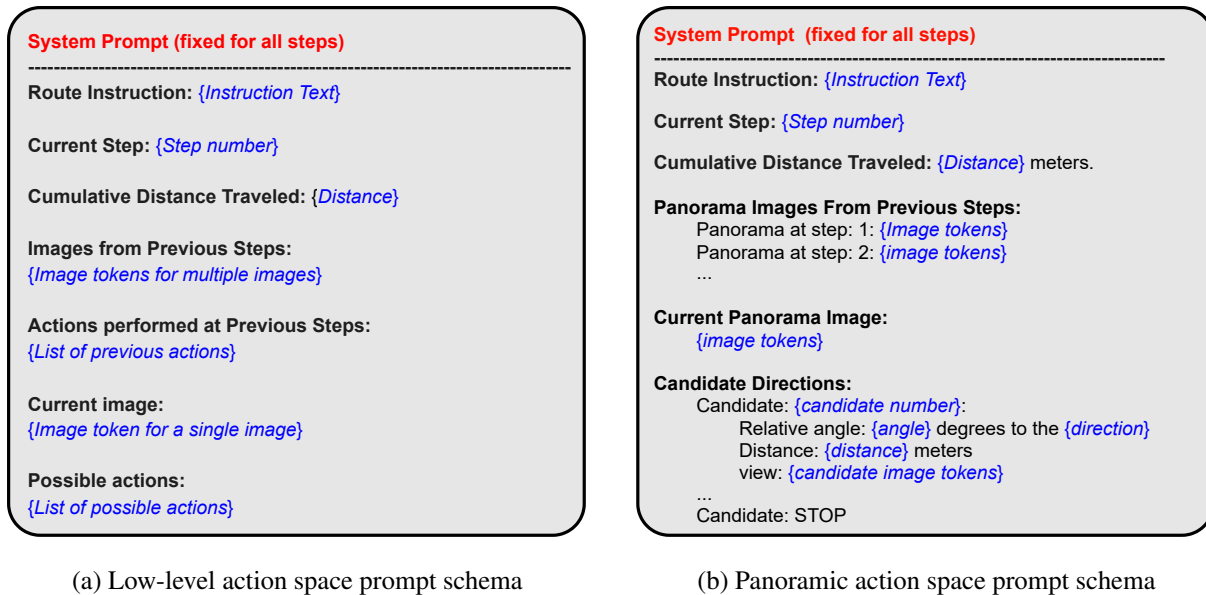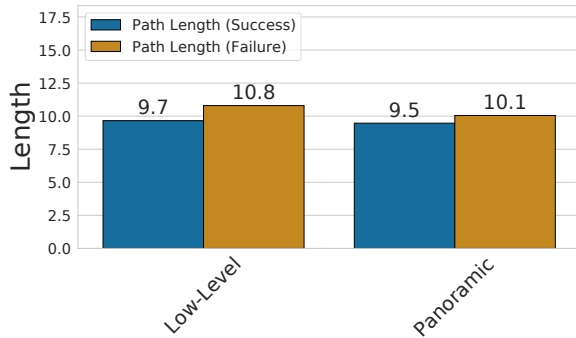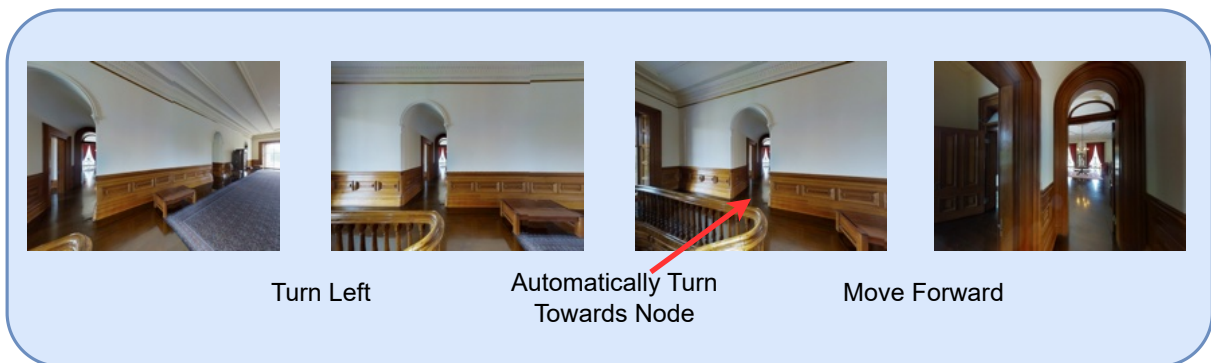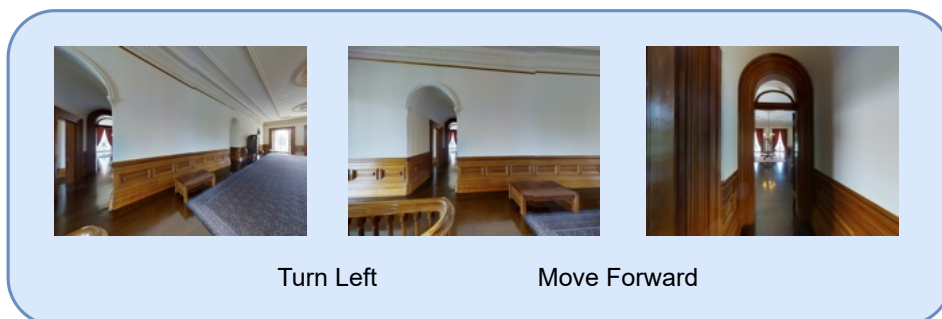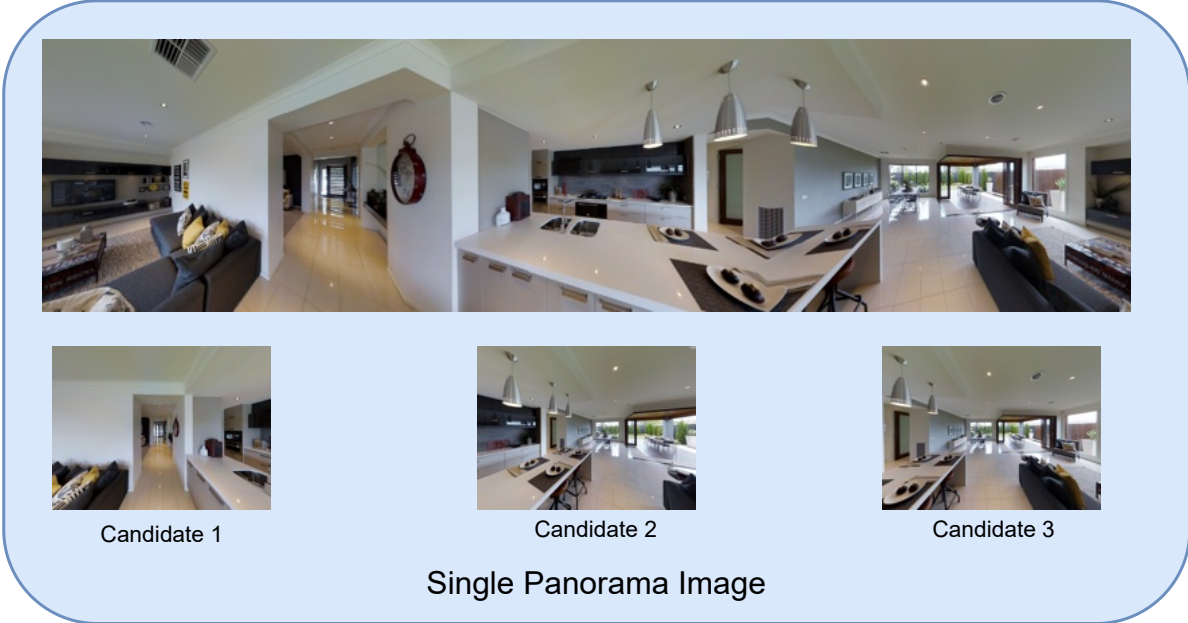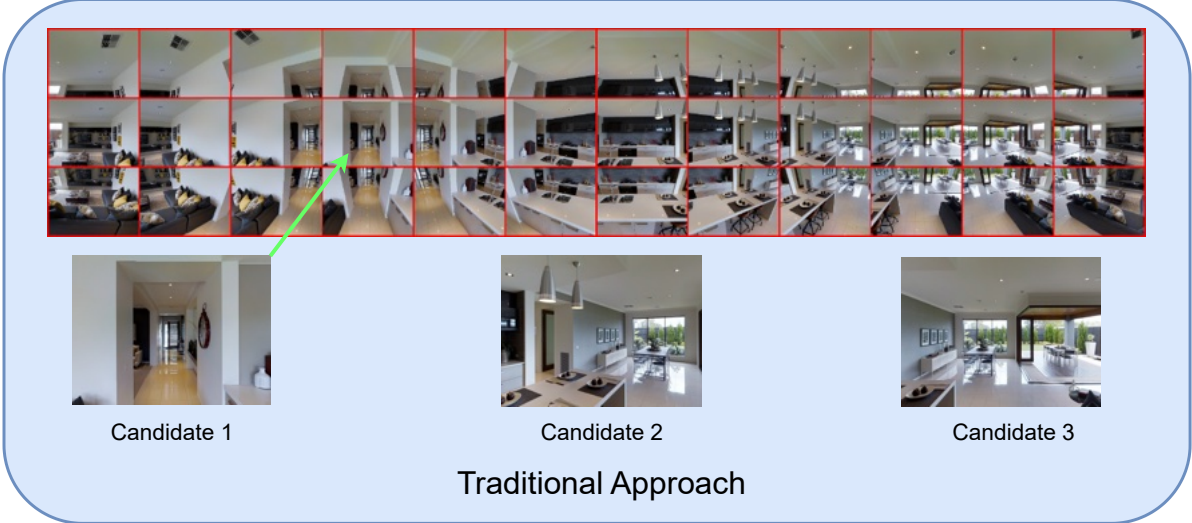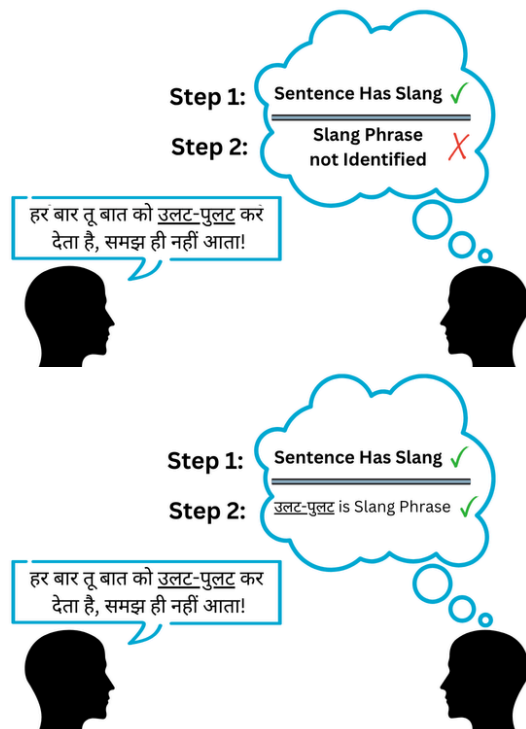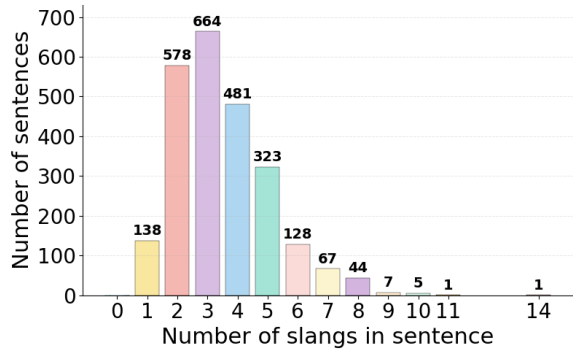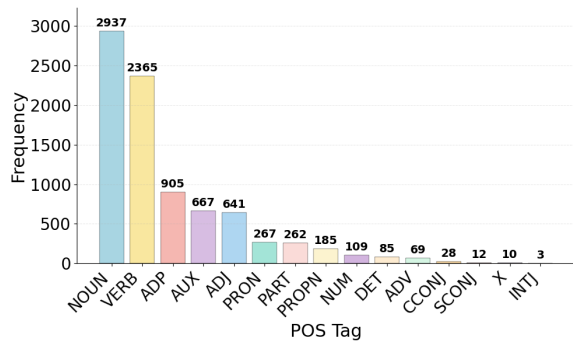{
    "multi_match": {
        "query": "query goes here",
        "fields": [
            "title^3",
            "abstract^2",
            "content"
        ],
        "type": "best_fields",
```

```
        "tie_breaker": 1
    }
}
```

**Filter Schema:** The following is the filter schema that is used to query the ElasticSearch index:

- **multi_match**: This is a query type in Elasticsearch that allows searching across multiple fields. It is particularly useful when the same query needs to be matched against various fields in the documents.

- **query**: Represents the search terms or phrases written by the user. This is the text that Elasticsearch will look for across the specified fields.

- **fields**:
  - **"title^3"**: This field corresponds to the full title of the document in Arabic. The ^3 denotes a boost factor of 3, meaning matches in the title are considered three times more relevant than matches in fields without a boost.
  - **"abstract^2"**: This field contains the abstract of the document in Arabic. With a boost factor of 2, matches here are deemed twice as relevant compared to unboosted fields.
  - **"content"**: This field encompasses the main content of the research documents, extracted from PDFs using either PyPDF or Surya. It does not have an explicit boost factor, so it serves as the baseline relevance.

- **type**: Set to **"best_fields"**, this parameter instructs Elasticsearch to return documents that have the best score from any one field. It focuses on the single most relevant field match

476

rather than combining scores from multiple fields.

- **tie_breaker**: With a value of **1**, the tie breaker adjusts the scoring when documents match the query in multiple fields. A higher tie breaker increases the influence of secondary matches on the overall score, ensuring that documents with multiple field matches are ranked higher.

More information about Elasticsearch Query Language can be found in the official documentation [8].

By leveraging this schema, we prioritized documents where the query terms appeared in the title or abstract, reflecting a higher likelihood of relevance. The boosting factors ensured that matches in the title and abstract had a more significant impact on the relevance score than matches in the main content. The use of the **best_fields** type, combined with an appropriate tie breaker, allowed for a balanced and effective retrieval of documents, enhancing the quality of the search results; after retrieving the search results from Elasticsearch, we extracted the relevance scores and the ranks of the documents as provided by the search engine. These scores and ranks are fed into the *relevation* tool for human evaluation.



Figure 1: System architecture of the proposed chatbot, illustrating the interaction between the IR, the LLM (OpenAI's GPT-4-turbo) and the auto-correct components for generating accurate, context-aware responses.

### 3.3 QA Evaluation

The classical way of evaluating QA systems requires gold dataset, which includes a set of questions and their answers. Then usually a metric

is based on lexical matching between the system-generated answer (prediction) and the real answer (reference) is used to generate a similarity score. Kamalloo et al. (2023) has shown that lexical-matching approaches suffer from a big fall, because of two reasons:

- It is almost impossible to provide a list of gold answers that cover all possibilities.

- Lexical matching methods cannot detect hallucinations and false claims in the output of LLM models.

Another work by Min et al. (2021) demonstrated that when humans who have experience in the domain evaluate the performance of question-answering (QA) system, the evaluation metric can increase up to $23\%$ compared to evaluations conducted through automated methods. Encouraged by their finding, we choose to perform human-evaluation of our proposed RAG system.

## 4  Experiments and Discussion

Figure 1 shows our system architecture. The system consists of mainly two components: the IR component and the Chatbot component. Users start their sessions by posing a question to the system, then the ChatBot helper generates queries from this question and conducts a search from the index. The purpose of the search is to build up context for the LLM model that will be used to generate a full answer to the user's question. This context helps reduce hallucination and makes inferring answers accurate by providing scientific references for users. The context consists of a list of papers (title, abstract and content) that will be injected to the system prompt. For that reason, a relevant context is highly important to generate answers that satisfy the user's question. To ensure high quality context, focus mainly on tuning choosing IR system that performs reasonably well for our use-case.

When indexing a new document, the following process is initiated: First, the document is parsed using an OCR system to extract its text. This text is then passed to the LLM with the proposed auto-correction prompt to fix any potential error introduced by the OCR. Finally, the corrected text is indexed by Elasticsearch, making it available for search queries.

---

[8]https://www.elastic.co/blog/getting-started-elasticsearch-query-language

## 4.1 IR experiments

To experiment with multiple scoring models in our dataset, we evaluate the performance of the IR system against the training dataset on the following retrieval model: LM-Dirichlet (Blei et al., 2003), Okapi BM25 (Crestani et al., 1998), and LM-JelinekMercer (Zhai and Lafferty, 2017).

A list of retrieval models supported by Elastic-Search is available in the official documentation[9]. We focus in our IR evaluation on two IR metrics: Normalized Discounted Cumulative Gain ($NDCG@k$), and Precision at k ($P@k$). $NDCG@k$ considers both the relevance and the position of the retrieved documents, with higher weights assigned to the results at the top of the list. $P@10$ measures the proportion of relevant items in the top k ranked documents, focusing solely on precision without considering the ranking order within the top-k items (Teufel, 2007).

Our IR experiments aim to address two pivotal questions:

- Can an off-the-shelf OCR model deliver satisfactory performance for Arabic IR, or is there a need for a more sophisticated model?

- Which IR model demonstrates superior performance in our experimental settings?

To address these two questions, we construct two ElasticSearch indices: one utilizing text extracted from PDF files using PyPDF library, and another employing the Surya model. The details of these indices are elaborated in Section 3.2. There is no difference between these two indices nor the querying mechanism.

Table 4 shows the evaluation results for the three IR models against the test set. This set is unseen and was not used to make any decision. The results show consistent improvement of both NDCG@5 and P@5 for all retrieval models when using more advanced OCR system (Surya), the biggest improvement is shown in Okapi BM25 model, where $NDCG@5$ increased by $+4.85$ points, and $P@5$ by $+7.53$ points. It is evident that Okapi BM25 demonstrates superior overall ranking quality ($NDCG@5$), LM-JelinekMercer shows slightly better precision in the top 5 results ($P@5$). This suggests that Okapi BM25 might

be more effective at distinguishing between degrees of relevance across the result set, while LM-JelinekMercer is particularly good at identifying highly relevant documents for the top positions, possibly due to its smoothing technique being well suited to the characteristics of Arabic scientific text. Considering those results, we decide to choose Okapi BM25 as the scoring model for our IR system, and Surya model as an OCR system to parse the PDF files.

## 4.2 Chatbot

In this section, we discuss how we use our IR system and leverage the LLM model to answer users' questions. First, user starts their session by posing a question to the system, as shown in Figure 1. Then the question is used as a query to retrieve top 5 relevant documents from Elasticsearch. Each document includes its title, abstract and content. These documents are then injected into the prompt, then a request is made to the Azure OpenAI API chat completion endpoint (*GPT-4 turbo-2024-04-09* version) [10] to answer the user's question. The LLM prompt is as follows:

> *You are an Expert Academic Research Synthesizer. Your function is to act as a research assistant, tasked with extracting and synthesizing information exclusively from a provided corpus of scientific documents. Strictly follow these instructions:*
> *1. **Corpus Definition**: You will receive a series of academic papers, each formatted with a clear title and content.*
> *2. **Information Adherence**: You MUST NOT use any external knowledge. All information in your response must be directly derived from the provided papers.*
> *3. **Answer Structure**: The response must be a comprehensive and cohesive synthesis of the information relevant to the question. Do not provide a list of facts; instead, integrate findings from multiple papers to create a single, detailed answer.*
> *4. **Specificity and Detail**: Focus on providing an extremely specific and factual answer. Avoid all forms of vague, generic, or abstract language.*

*5. **Citations**: Every single statement of fact or claim must be followed immediately by a citation in the format '(Paper-Title)'. If a sentence synthesizes information from multiple papers, list all relevant citations.*

*6. **Language**: The final, complete answer MUST be generated entirely in the Arabic language.*

*7. Provide only the answer to the question. Given the following papers:*
*{{Title   Abstract   Content}}*

*.*

*.*

*{{Title Abstract    Content}}*
*Answer the following question: {{Question}}*

The prompt is designed to make the LLM simulate the process of a human researcher while strictly preventing hallucination. By forbidding external knowledge and demanding that every factual statement be anchored to a source with a mandatory citation, we enforce a high degree of verifiability, ensuring the model cannot invent information. This mimics a researcher's reliance on primary sources. Moreover, the instruction to integrate findings from multiple papers to create a single, detailed answer compels the model to move beyond simple fact extraction and replicate the human cognitive process of synthesis. This dual approach ensures that the generated output is not only factually grounded and trustworthy but also demonstrates a sophisticated, human-like understanding of the source material, a crucial requirement for reliable academic use.

To evaluate the performance of our proposed Chatbot architecture, we developed two Telegram bots utilizing the Telegram Bot API [11]. This API allows for the creation of programs that use Telegram messages as an interface. Users can interact with it using their mobile devices or Telegram desktop version. Our experimental setup consisted of two distinct bots:

- **Baseline Bot** This bot directly sends user questions to the LLM endpoint. The only system prompt that we use in the baseline system is: *Answer the following question and provide citations for your answer*: $\{question\}$.

- **The Proposed Chatbot (ArabicRAG)** This bot implements our proposed architecture (as

described in Figure 1), utilizing the prompt we presented earlier, including contextual information.

Using the Telegram API, we were able to create a robust experimental framework to compare the performance of our proposed architecture with the baseline system. We asked our human judges to send each question in the test set to both systems and give each answer a grade based on its relevance. We introduce the following grading system, the final grade is the sum of all of them:

- **Does the generated output provide correct citations?**
    - 3: There is at least one correct citation for each statement.
    - 2: Some correct citations are missing, but not very crucial.
    - 1: Crucial citations are missing, or incorrect citations are provided.
    - 0: There are no correct citations provided.

- **Does the generated output answer your question?**
    - 3: Yes, the output fully answers my question.
    - 2: The output partially answers my question.
    - 1: The output somewhat answers my question.
    - 0: The output does not answer my question.

Then we take the average of the two grades, a perfect answer will be graded 3 for citation and 3 for correctness, yielding an average of 3 final grade. This evaluation is done for the 70 queries in the test set.

Table 5 summarizes the performance of our proposed system (ArabicRAG) and the baseline (GPT-4o). ArabicRAG demonstrates superior performance across all the metrics. With a 60% higher rate of fully correct answers (score=3 as judged by the human experts), and 3 times fewer complete failures (score=0) compared to GPT. We can notice a smaller standard deviation in ArabicRAG (0.8) compared to GPT of 1.2. This means that ArabicRAG tends to have more predictable performance around the mean score. GPT in 15% of the test

| Model | PyPDF OCR | | Surya OCR | |
|---|---|---|---|---|
| | NDCG@5 | P@5 | NDCG@5 | P@5 |
| Okapi BM25 | 67.30 | 37.97 | **73.15** | **45.50** |
| LMDirichlet | 58.99 | 38.82 | 61.60 | 41.39 |
| LM-JelinekMercer | 65.84 | 41.55 | 69.17 | 43.10 |

Table 4: Performance of three retrieval models against the test dataset, both NDCG@10 and P@10 are reported in percentages, bold numbers are statistically significants

| Metric | ArabicRAG | Baseline |
|---|---|---|
| Mean Score ($\pm SD$) | 2.4 ($\pm 0.8$) | 1.5 ($\pm 1.2$) |
| %Score =3 | 60% | 40% |
| %Score =0 | 5% | 15% |

Table 5: Evaluation of the two bots on the test set. We report average of relevance graded score, and average of citation graded score for each system

set (around 10 questions) fails to provide a correct answer with accurate citations. For example, when GPT is asked to provide a brief introduction about historical figures from the Arabic literature, it tends to provide basic information as usually found in Wikipedia, but the main issue is with almost unrelated citations to books and articles. When we checked those references, we found out that they are irrelevant. This is considered hallucination. The pattern of synthetic scholarships poses particular risks in academic applications where source authenticity is crucial to the credibility of the research.

However, upon examining the failures in the ArabicRAG system, we observed that the inability to generate correct answers is primarily due to the limited number of research documents in the corpus (approximately 10,000). Consequently, the Information Retrieval system often retrieves documents that are poorly relevant to the questions, leading to irrelevant answers generated by the LLM. One solution is to enable real-time internet searches from reliable Arabic sources when the retrieved documents have low similarity scores. Another potential solution is to continuously expand the corpus by indexing more documents, thereby covering a diverse set of research topics. This is a potential future work of this research.

## 5 Conclusion

In this paper, we explored the impact of advanced Arabic language preprocessing techniques on the performance of information retrieval systems and

their downstream influence on retrieval-augmented generation systems. Our findings suggest that employing a state-of-the-art deep learning-based OCR system (Surya) significantly enhances the IR performance, with improvements of up to 8 points in $P@5$ and 11% in RAG answering accuracy compared to baseline system. These results underscore the importance of robust preprocessing and language-aware IR in addressing challenges posed by Arabic script complexity and domain-specific terminology.

By integrating our enhanced IR system with a generative model, we developed ArabicRAG, a chatbot capable of providing contextually accurate and citation-rich answers to academic queries. Comparative evaluations against a baseline system revealed that ArabicRAG achieves a 20% higher rate of fully correct answers and significantly reduces hallucinations.

Future work will focus on expanding the corpus to cover a broader range of research topics and exploring real-time internet-based retrieval to address low-similarity cases. These enhancements aim to further improve the system's ability to deliver accurate and relevant responses, thereby advancing the state of Arabic NLP in academic contexts.

## References

Gagan Bhatia, El Moatez Billah Nagoudi, Fakhraddin Alwajih, and Muhammad Abdul-Mageed. 2024. Qalam: A multimodal LLM for Arabic optical character and handwriting recognition. In *Proceedings of The Second Arabic Natural Language Processing Conference*, pages 210–224, Bangkok, Thailand. Association for Computational Linguistics.

David M Blei, Andrew Y Ng, and Michael I Jordan. 2003. Latent dirichlet allocation. *Journal of machine Learning research*, 3(Jan):993–1022.

Tom B. Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, Sandhini Agarwal, Ariel Herbert-Voss, Gretchen Krueger, Tom Henighan, Rewon Child,

Aditya Ramesh, Daniel M. Ziegler, Jeffrey Wu, Clemens Winter, Christopher Hesse, Mark Chen, Eric Sigler, Mateusz Litwin, Scott Gray, Benjamin Chess, Jack Clark, Christopher Berner, Sam McCandlish, Alec Radford, Ilya Sutskever, and Dario Amodei. 2020. Language models are few-shot learners. In *Proceedings of the 34th International Conference on Neural Information Processing Systems*, NIPS '20, Red Hook, NY, USA. Curran Associates Inc.

H. Bunke, M. Roth, and E.G. Schukat-Talamazzini. 1995. Off-line cursive handwriting recognition using hidden markov models. *Pattern Recognition*, 28(9):1399–1413.

Fabio Crestani, Mounia Lalmas, Cornelis J. Van Rijsbergen, and Iain Campbell. 1998. "is this document relevant?... probably": a survey of probabilistic models in information retrieval. *ACM Comput. Surv.*, 30(4):528–552.

Dima El Zein and Célia da Costa Pereira. 2022. User's knowledge and information needs in information retrieval evaluation. In *Proceedings of the 30th ACM Conference on User Modeling, Adaptation and Personalization*, UMAP '22, page 170–178, New York, NY, USA. Association for Computing Machinery.

Ehsan Kamalloo, Nouha Dziri, Charles Clarke, and Davood Rafiei. 2023. Evaluating open-domain question answering in the era of large language models. In *Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 5591–5606, Toronto, Canada. Association for Computational Linguistics.

Tom Kocmi and Christian Federmann. 2023. Large language models are state-of-the-art evaluators of translation quality. In *Proceedings of the 24th Annual Conference of the European Association for Machine Translation*, pages 193–203, Tampere, Finland. European Association for Machine Translation.

Sewon Min, Jordan Boyd-Graber, Chris Alberti, Danqi Chen, Eunsol Choi, Michael Collins, Kelvin Guu, Hannaneh Hajishirzi, Kenton Lee, Jennimaria Palomaki, et al. 2021. Neurips 2020 efficientqa competition: Systems, analyses and lessons learned. In *NeurIPS 2020 Competition and Demonstration Track*, pages 86–111. PMLR.

Fabio Petroni, Federico Siciliano, Fabrizio Silvestri, and Giovanni Trappolini. 2024. IR-RAG @ SIGIR24: Information Retrieval's Role in RAG Systems. In *Proceedings of the 47th International ACM SIGIR Conference on Research and Development in Information Retrieval*, SIGIR '24, page 3036–3039, New York, NY, USA. ACM.

Jawad Sadek, Fairouz Chakkour, and Farid Meziane. 2012. Arabic Rhetorical Relations Extraction For Answering "Why" and "How to" Questions. In *Proceedings of the 17th International Conference on Applications of Natural Language Processing and Information Systems*, NLDB'12, page 385–390, Berlin, Heidelberg. Springer.

Simone Teufel. 2007. *An Overview of Evaluation Methods in TREC Ad Hoc Information Retrieval and TREC Question Answering*, pages 163–186. Springer Netherlands, Dordrecht.

Chengxiang Zhai and John Lafferty. 2017. A study of smoothing methods for language models applied to ad hoc information retrieval. In *ACM SIGIR Forum*, volume 51, pages 268–276. ACM New York, NY, USA.

# 8th International Conference on Natural Language and Speech Processing

August 25-27, 2025 Odense, Denmark

# Technical program

## Day 1, Monday 25 August

| | |
|---|---|
| **08.30 - 09.15** | **Registration & Breakfast** |

| | |
|---|---|
| **09.15 - 10.30** | **ICNLSP 2025 Opening**<br>**Keynote by Prof. Dr. Barbara Plank**<br>*Human-Centered LLMs for Inclusive Language Technology* |

| | |
|---|---|
| **10.30 - 11.00** | **Coffee Break** |

| | |
|---|---|
| **11:00 - 13:00** | **Session 1 (Chair: Torben Johansen)** |

- *L1RA: Dynamic Rank Assignment in LoRA Fine-Tuning*
  Raul Singh; **Nicolò Brunello**; Vincenzo Scotti; Mark Carman

- *A Comparison between Low-level and Panoramic Action Spaces*
  **Vebjørn Kåsene**; Pierre Lison

- *The Impact of Annotator Personas on LLM Behavior Across the Perspectivism Spectrum*
  **Olufunke Sarumi**; Charles Welch; Daniel Braun; Jörg Schlötterer

- *Efficient Continual Learning for Small Language Models with a Discrete Key-Value Bottleneck*
  **Andor Diera**; Lukas Galke; Fabian Karl; Ansgar Scherp

- *Word-level Language Identification using Encoder-Only and Decoder-Only Models*
  Javier Iranzo-Sánchez; Parnia Bahar; Alejandro Pérez-González-de-Martos; **Mattia Antonino Di Gangi**

| | |
|---|---|
| **13:00 - 14:00** | **Lunch Break** |

## 14:00 - 15:30    Session 2 (Chair: Lukas Galke)

- *Improving French Synthetic Speech Quality via SSML Prosody Control*
  **Nassima Ould Ouali**; Awais Hussain Sani; Ruben Bueno; Jonah Dauvet; Tim Luka Horstmann; Eric Moulines

- *Adapting ASR Models for Speech-to-Punctuated-Text Recognition with Utterance Gluing*
  **Agata Jakubiak**; Piotr Stachyra; Piotr Czubowski; Hubert Borkowski; Sebastian Łątka ; Radosław Iżak; Kornel Jankowski; Sonia Janicka; Mateusz Zieliński

- *Breaking the HISCO Barrier*
  **Torben Johansen**; Christian Vedel; Christian Dahl

- *CMC-SC: Cross-Modal Contextualized ASR Spelling Correction via BERT and WavLM using a Soft Fusion Framework*
  Mohammad Reza Peyghan; Sajjad  Amini; **Shahrokh Ghaemmaghami**

- *Bridging the Gap: Design and Evaluation of an Automated System for French Cued Speech*
  **Brigitte Bigi** (CNRS)

## 15:30 - 16:00    Coffee Break

## 16:00 - 17:40    Session 3 (Chair: Jakob Kusnick)

- *Zero-Shot Commonsense Validation and Reasoning with Large Language  Models: An Evaluation on SemEval-2020 Task 4 Dataset*
  **Mohammad AL-Smad**i; Rawand  Alfugaha

- *The Chunking Paradigm: Recursive Semantic for RAG Optimization*
  **Seemab Latif**; Huma Ameer; Hannan Akram; Mehwish Fatima

- *Efficient Long-Document Summarization with Mixture of LoRA Experts, Sparse Attention and Multi-Scale Rotary Positional Embeddings*
  **P. Yadla**

- *From Context to Emotion: Leveraging LLMs for Recognizing Implicit Emotions*
  **Hanane Boutouta**; Abdelaziz Lakhfif; Ferial Senator ; Chahrazed Mediani

- *From Outliers to Topics in Language Models: Anticipating Trends in News Corpora*
  **Evangelia ZVE**; Benjamin Icard; Alice Breton; Lila Sainero; Gauvain Bourgne ; Jean-Gabriel Ganascia

## 17:40 - 20:00    Posters Session + Reception

- *From Performance to Process: Temporal Information Dynamics in Language Model Fine-tuning*
  **Frida Hæstrup**; Ross Deans Kristensen-McLachlan

- *Next Speaker Prediction for Multi-Speaker Dialogue with Large Language Models*
  **Lukas Hilgert**; Jan Niehues

- *Beyond Shallow Heuristics: Leveraging Human Intuition for Curriculum Learning*
  **Vanessa Toborek***; Sebastian Müller; Tim Selbach; Tamás Horváth; Christian Bauckhage

- *Speech-Based Depressive Mood Detection in the Presence of Multiple Sclerosis*
  **Monica Gonzalez Machorro**; Uwe Reichel; Pascal Hecker; Helly Hammer; Hesam Sagha; Florian Eyben; Robert Hoepner; Björn Schuller

- *topicwizard - a Modern, Model-agnostic Framework for Topic Model Visualization and Interpretation*
  **Márton Kardos**; Kenneth Enevoldsen; Kristoffer Nielbo

- *Tokenization and Morphology in Multilingual Language Models: A Comparative Analysis of mT5 and ByT5*
  Thao Anh Dang; Limor Raviv; **Lukas Galke**

- *Mapping Educational Science in the Broader Academic Discourse about Artificial Intelligence*
  Martin Rehm; **Zhiru Sun**; Maxime Holmberg Sainte-Marie

- *Re-Representation in Sentential Relation Extraction with Sequence Routing Algorithm*
  **Ramazan Bahrami**; Ramin Yahyapour

## Day 2, Tuesday 26 August

**08.30 - 09.00   Registration & Breakfast**

**09.00 - 10.40   Session 4 (Chair: Tariq Yousef)**

- *Domain adaptation and question-answer pooling for Aphasia modelling*
  Uwe Reichel; **Monica Gonzalez Machorro**; Lisa Maria Ehlen; Pascal Hecker; Dorothea Peitz; Cornelius Werner; Felix Burkhardt; Christian Kohlschein; Florian Eyben; Björn Schuller

- *GUIDE: A Framework for Improving Functional Software Test Descriptions with Language Models*
  **Mathis Ronzon**; Thierry Roger; Zoltan Miklos; Annie Foret

- *Aligning University Curriculum with Market Needs: A Longitudinal NLP Study of Danish Academic Skill Demands*
  **Zhiru Sun**; Jacob Mørup Wang

- *Towards Robust Urdu Aspect-based Sentiment Analysis through Weakly-Supervised Annotation Framework*
  Zoya Maqsood; **Seemab Latif**; Rabia Latif

**10.40 - 11.00   Coffee Break**

**11:00 - 12:00   Keynote by Prof. Dr. Anders Søgaard**
**What to think of NLP these days?**

**12:00 - 13:00   Lunch Break**

**13:00 - 14:40   Session 5 (Chair: Zhiru Sun)**

- *FActBench: A Benchmark for Fine-grained Automatic Evaluation of LLM-Generated Text in the Medical Domain*
  **Anum Afza**l; Juraj Vladika; Florian Matthes

- *Cross-Lingual Sentence-Level Skill Identification in English and Danish Job Advertisements*
  **Nurlan Musazade**; Mike Zhang; József Mezei

- *[LLMs Information Flow Diagnostic: Memory-based Evidence from Random Matrix Theory](#)*
  **Sami Diaf**

- *[A Retail-Corpus for Aspect-Based Sentiment Analysis with Large Language Models](#)*
  **Oleg Şilcenco**; Marcos Machado; Wallace Ugulino; Daniel Braun

- *[Demographics and Democracy: Benchmarking LLMs' Gender Bias and Political Leaning in European Parliament](#)*
  **Jinrui Yang**; Xudong Han ; Timothy Baldwin

### 14:40 - 15:00    Coffee Break

### 15:00 - 17:00    Session 6 (Chair: Esben Andreas Wrona Bay Sørensen)

- *[ASR Models for Traditional Emirati Arabic: Challenges, Adaptations, and Performance Evaluation](#)*
  **Maha Alblooki**

- *[Building an Ewe Language Dataset: Towards Enhancing Automatic Speech Recognition Technologies for Low Resource Languages](#)*
  Isaac Wiafe; Akon Obu Ekpezu; **Raynard Dodzi Helegah**; Fiifi Baffoe Payin Winful ; Elikem Doe Atsakpo; Charles Nutrokpor; Kafui Kwashie Solaga

- *[CLEAR: Code-Mixed ASR with LLM-Driven Rescoring](#)*
  **Shivam Kumar**; Md Shad Akhtar

- *[Evaluating ASR in a Clinical Context : What Whisper Misses](#)*
  Haeeul Hwang; **Eric Jordan**; Deok-Hee Kim-Dufor; Christophe Lemey; Motasem Alrahabi

- *[Assessing ASR Robustness for Burmese: Impacts of Missing Speech Segments and Interruptions](#)*
  **Ankit Maurya**; Manikandan Ravikiran; Rohit Saluja

### 18:00        Reception at Storm's Pakhus

**Address:** Lerchesgade 4, 5000 Odense, Denmark

## Day 3, Wednesday 27 August

| | |
|---|---|
| **08.30 - 09.00** | **Registration & Breakfast** |

| | |
|---|---|
| **09.00 - 10.40** | **Session 7 (Chair: Alexandra Diehl)** |

- *Scalable Text Vectorization with Hyperdimensional Computing Through Selective Word Encoding*
  **Timur Mudarisov**; Evgeny Polyachenko; Zsofia Kraussl; Enriqueta Patricia Becerra Sanchez; Tatiana Petrova; Radu State

- *Style-Controlled Response Generation for Dialog Systems with Intimacy Interpretation*
  **Takuto Miura**; Kiyoaki Shirai; Natthawut Kertkeidkachorn

- *Dora explores Clinically Relevant Information in EHRs using NER*
  **Martin Laursen**; Lina Elkjær Pedersen; Josefine Bak H. Adelhelm; Rasmus Bank Lynggaard; Pernille Just Vinholt

- *Beyond Labeled Datasets: Advancing TTS with Direct Preference Optimization on Unlabeled Speech Dataset*
  **Andrii Zhuravlov**; Volodymyr Sydorskyi

- *CUPE: Contextless Universal Phoneme Encoder for Language-Agnostic Speech Processing*
  Abdul Rehman; Jian-jun Zhang; Xiasong Yang

| | |
|---|---|
| **10.40 - 11.00** | **Coffee Break** |

| | |
|---|---|
| **11:00 - 12:00** | **Keynote by Prof. Dr. Peter Schneider-Kamp** <br> *The Cost of Intelligence: Efficiency Is the Only Path to Democratized AI* |

| | |
|---|---|
| **12:00 - 13:00** | **Lunch Break** |

## 13:00 - 15:00  Session 8 (Chair: Nicklas Sindlev Andersen)

- *Tachelhiyt-Darija: a parallel speech corpus for two underrepresented languages*
  Noureddine Atouf; **Elsayed Issa**; Said Ouzbayr

- *The Need for Robust and Inclusive Benchmarks in Evaluating LLMs on Arabic Text*
  **Lubana Al Rayes**; Ashraf Elnagar

- *On Limitations of LLM as Annotator for Low Resource Languages*
  Raviraj Joshi; Abhay Shanbhag; Amogh Thakurdesai; Ridhima Sinare; **Suramya Jadhav**

- *HiSlang-4.9k: A Benchmark Dataset for Hindi Slang Detection and Identification*
  **Tanmay Tiwari**; Vibhu Gupta; Manikandan Ravikiran; Rohit Saluja

- *Enhancing Arabic Retrieval Augmented Generation through Language Processing*
  **Shadi Saleh**; Belkacem EL Jattari; Layth Oud; Maryam Alblooshi; Bouthaina Lakhdari; Ali Alnaqbi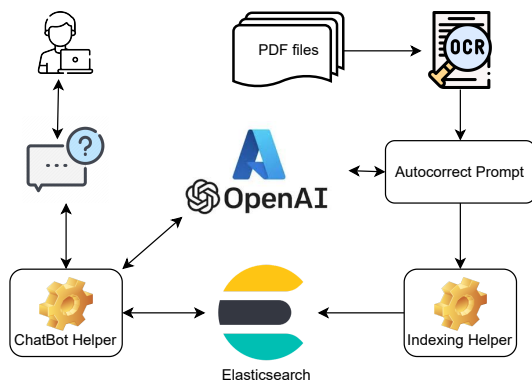